

# Parameter estimation in a subcritical percolation model with colouring

Felix Beck<sup>\*†</sup>, Bence Mélykúti<sup>\*‡</sup>

13th March 2017

## Abstract

In the bond percolation model on a lattice, we colour vertices with  $n_c$  colours independently at random according to Bernoulli distributions. A vertex can receive multiple colours and each of these colours is individually observable. The colours colour the entire component into which they fall. Our goal is to estimate the  $n_c + 1$  parameters of the model: the probabilities of colouring of single vertices and the probability with which an edge is open. The input data is the configuration of colours once the complete components have been coloured, without the information which vertices were originally coloured or which edges are open.

We use a Monte Carlo method, the method of simulated moments to achieve this goal. We prove that this method is a strongly consistent estimator by proving a uniform strong law of large numbers for the vertices' weakly dependent colour values. We evaluate the method in computer tests. The motivating application is cross-contamination rate estimation for digital PCR in lab-on-a-chip microfluidic devices.

**Keywords** parameter estimation, method of simulated moments, percolation, strong law of large numbers with dependence, microfluidics, cross-contamination

**Mathematics subject classification** 62F10 (Point estimation), 60K35 (Interacting random processes; statistical mechanics type models; percolation theory)

## 1 Bond percolation with colouring

We consider bond percolation [14] on the triangular lattice, but our arguments hold for the square lattice as well. The vertex set of the infinite lattice is denoted by  $L$ . Edges are open (that is, included in the graph, alternatively, receive weight 1 as opposed to 0) independently at random with probability  $\mu \in [0, 1]$ . There are  $n_c \in \mathbb{N} \setminus \{0\}$  colours given, and for every colour  $\ell \in \{1, 2, \dots, n_c\}$ , a parameter  $\lambda^\ell \in [0, 1]$  is fixed. For every vertex  $i \in L$ , the vertex is coloured with colour  $\ell \in \{1, 2, \dots, n_c\}$  according to a Bernoulli random variable with probability  $\lambda^\ell$ . The colouring with different colours is independent in any one vertex, and it is also independent among different vertices. A vertex can receive multiple colours and each of these colours is individually observable. We call this colouring the *seeding*:  $X_i^\ell \in \{0, 1\}$  for every  $i \in L$  and  $\ell \in \{1, 2, \dots, n_c\}$ .

These colours propagate through open edges and colour ('contaminate') the entire component they are contained in. Let  $i \leftrightarrow j$  mean that vertices  $i, j \in L$  are connected by an open path. The

---

<sup>\*</sup>Centre for Biological Systems Analysis (ZBSA), University of Freiburg, Habsburgerstraße 49, 79104 Freiburg, Germany

<sup>†</sup>Institute for Mathematics, University of Freiburg, Eckerstraße 1, 79104 Freiburg, Germany

<sup>‡</sup>Corresponding author. Email: [melykuti@stochastik.uni-freiburg.de](mailto:melykuti@stochastik.uni-freiburg.de).

observed colour configuration is

$$Y_i^\ell := X_i^\ell \vee \bigvee_{\substack{j \in L \\ j \leftrightarrow i}} X_j^\ell \in \{0, 1\}$$

for every  $i \in L$  and  $\ell \in \{1, 2, \dots, n_c\}$ , where  $\vee$  is the maximum operator.

We also consider this process on finite, connected subsets of the lattice,  $I \subset L$ . (Here connected is meant with all lattice edges considered, not only the open edges.) Picking the vertex set  $I$  implicitly fixes its edge set, the edges which connect vertices of  $I$ . We let  $n_I := |I|$ . We write  $i \sim j$  for adjacent lattice vertices  $i, j \in L$  no matter in what state the connecting edge is.

Often we consider nested sequences of such  $I$  where each successor is a superset of its predecessor and  $n_I \rightarrow \infty$ . We fix an ordering of the vertices of the infinite lattice  $L$  which is compatible with this sequence as  $n_I \rightarrow \infty$ , that is, each  $I$  comprises vertices labelled with  $\{1, \dots, n_I\}$ . We use  $I_2 := \{(i, j) \in I \times I \mid i \sim j, i < j\}$  for the set of ordered pairs of adjacent vertices (independently of whether the connecting edge is open or closed) and  $n_p := |I_2|$  for the total number of possible edges within  $I$ . We define the *exterior vertex boundary* of a subset  $I$  by

$$\Delta I := \{j \in L \mid j \notin I, \exists i \in I : i \sim j\}.$$

We always require that in our sequences,  $|\Delta I|/|I| \rightarrow 0$  and for the triangular lattice,  $n_p \sim 3n_I$  (asymptotic equality;  $n_p \sim 2n_I$  is the corresponding condition for the square lattice).

For a fixed  $I$ , we define a variant of  $Y_i^\ell$  that is determined exclusively by the seeding and edges in  $I$ :

$$\tilde{Y}_i^\ell := X_i^\ell \vee \bigvee_{\substack{j \in I \\ j \leftrightarrow i}} X_j^\ell \in \{0, 1\}$$

for every  $i \in I$  and  $\ell \in \{1, 2, \dots, n_c\}$ . Here  $\leftrightarrow$  means connectedness by open edges in the edge set of  $I$ .

Our goal is to estimate the parameter  $\theta = (\lambda^1, \dots, \lambda^{n_c}, \mu)$  from the data  $(\tilde{Y}_i^\ell)_{i \in I, \ell \in \{1, 2, \dots, n_c\}}$  (Figure 1). The spatial arrangement of  $(\tilde{Y}_i^\ell)$  within the lattice is known, but the seeding  $(X_i^\ell)$  and the open or closed state of the edges are unavailable. We bring together four theoretical tools in this paper.

First, parameter estimation is conducted by the *method of simulated moments* (MSM) [12, 13] (Section 2). This is a simulation-based, computationally intensive statistical method that yields a point estimate for  $\theta$  which converges almost surely to the correct value as  $n_I \rightarrow \infty$ .

Second, as the first step towards proving the strong consistence of the estimator, we prove a strong law of large numbers (SLLN) with weakly dependent variables. We do this in Section 3 by adapting Theorem 1 of [8].

Third, the SLLN result requires some grasp of how small the dependence is between distant vertices of the lattice. The upper bounds on correlations are provided by the FKG and BK inequalities of percolation theory and the exponential decay of the cluster size distribution [1, 4, 9], [14, Chapters 2 and 6] in Section 4.

Fourth, for the strong consistence of the estimator, we extend the SLLN to be uniform in the parameter vector. We verify in Section 5 that the conditions of a sufficient condition for the uniform law of large numbers (ULLN) hold [19, p. 8, 2 Theorem] [22, p. 25, Lemma 3.1].

Our estimation method is tested on synthetic data with known parameter values in Section 6 and its performance is evaluated. In Section 7, the motivating problem is described, and the paper concludes with a discussion of possible improvements in modelling and methodology.

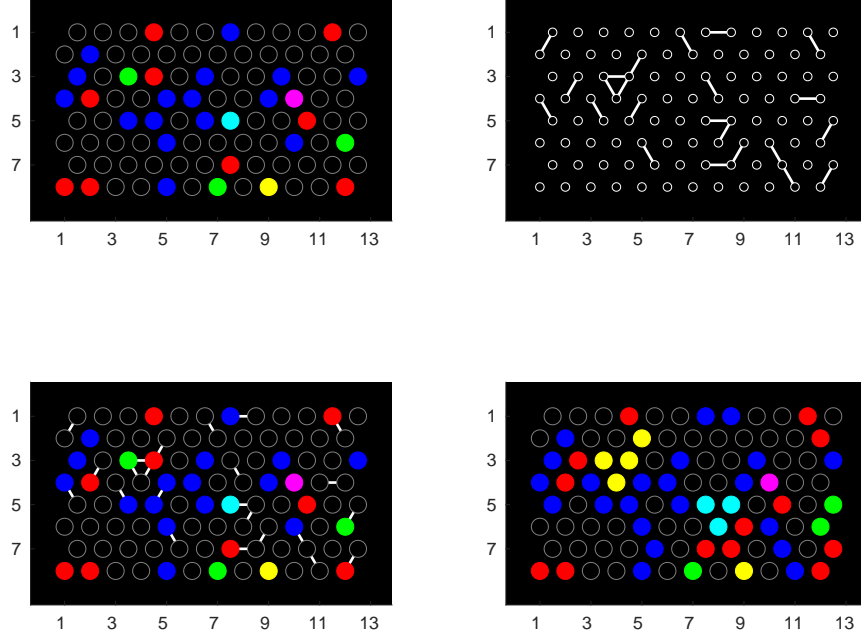


Figure 1: (*top left*) A realisation of random seeding ( $X_i^\ell$ ) with  $(\lambda^{\text{red}}, \lambda^{\text{green}}, \lambda^{\text{blue}}) = (0.1, 0.05, 0.2)$ . (*top right*) A realisation of bond percolation on the triangular lattice with  $\mu = 0.1$ . (*bottom left*) The bond percolation overlaid with the seeding. (*bottom right*) The resulting configuration ( $\tilde{Y}_i^\ell$ ) which serves as the data.

## 2 Method of simulated moments (MSM)

The MSM is a modification of the classical method of moments for parameter estimation for the case when the moments of the sampling distribution cannot be computed from the parameters in closed form. The MSM proposes to simulate  $n_s$  independent, identically distributed samples from the distribution, repeatedly with different parameter values  $\theta$  (usually, but not strictly necessarily, with common random variables as  $\theta$  is changed), and to choose the  $\theta$  which gives the closest match between moments of the data and that of the simulated data. For its detailed description, we recommend perusing a combination of [13] and [12].

The data  $\mathcal{Y} = (\mathcal{Y}_i)_{i \in I}$  originates from a distribution which is parameterised by the unknown  $\theta_0 \in \Theta$ .  $\theta_0$  is called the true value of the parameter. Normally, the  $\mathcal{Y}_i$  are independent. A sample from this family of distributions with a general parameter is denoted by  $Y = (Y_i)_{i \in I}$ . Let  $K$  be some  $n_m$ -dimensional function of the individual observations  $Y_i$ . Let  $k(\theta)$  be the expectation of  $K$  when  $K$  is evaluated on a draw  $Y_i$  from the distribution with parameter  $\theta \in \Theta$ ,  $k(\theta) := \mathbb{E}_\theta[K(Y_i)]$ . Thus  $k$  is a vector of  $n_m$  generalised moments of the distribution of  $Y_i$ . ( $\mathbb{E}_\theta$  is the expectation under the distribution with parameter  $\theta$ . Similarly,  $P_\theta$  is the probability of an event in that case.)

Let  $g$  be some multidimensional function that represents estimating constraints. In our case these are distances between observed moments and moments of the model with given parameter value  $\theta$ :

$$g(\mathcal{Y}_i, \theta) = K(\mathcal{Y}_i) - k(\theta).$$

By introducing  $\mathbb{E}_0$  as a shorthand for  $\mathbb{E}_{\theta_0}$ , it is immediate that  $\mathbb{E}_0[g(\mathcal{Y}_i, \theta_0)] = 0$ . However, for the

parameter estimation problem to be well posed, we require that

$$\mathbb{E}_0[g(\mathcal{Y}_i, \theta)] = 0 \quad \Longleftrightarrow \quad \theta = \theta_0. \quad (1)$$

Implicit in this is that we have at least as many independent equations as parameters.

The MSM is used when  $k(\theta)$  is not available in analytical form but there exists an unbiased estimator  $\tilde{k}(U_i^s, \theta)$ , and consequently an unbiased estimator for  $g$ ,  $\tilde{g}(\mathcal{Y}_i, U_i^s, \theta) = K(\mathcal{Y}_i) - \tilde{k}(U_i^s, \theta)$ . Here  $(U_i^s)_{i \in I, s \in \{1, \dots, n_s\}}$  is some source of randomness, typically vectors of independent, uniform random variables on  $[0, 1]$  as provided by a pseudorandom number generator. The estimators satisfy  $\mathbb{E}[\tilde{k}(U_i^s, \theta)] = k(\theta)$  and  $\mathbb{E}[\tilde{g}(\mathcal{Y}_i, U_i^s, \theta) \mid \mathcal{Y}_i] = g(\mathcal{Y}_i, \theta)$ .

We introduce a weighting by a symmetric, positive definite matrix  $\Omega \in \mathbb{R}^{n_m \times n_m}$ , which might be a function of the data, and consider the quadratic form  $\alpha(\eta) = \eta^T \Omega \eta$ . The broad principle of the MSM is the following.

**Proposition 1** The MSM estimator is defined as

$$\hat{\theta}_{n_s, n_I} := \arg \min_{\theta \in \Theta} \alpha \left( \frac{1}{n_I} \sum_{i=1}^{n_I} \left( K(\mathcal{Y}_i) - \frac{1}{n_s} \sum_{s=1}^{n_s} \tilde{k}(U_i^s, \theta) \right) \right).$$

If identifiability holds,  $n_s$  is fixed and  $n_I$  tends to infinity, and the almost sure convergence guaranteed by the SLLN

$$\frac{1}{n_I} \sum_{i=1}^{n_I} \tilde{k}(U_i^s, \theta) \xrightarrow{n_I \rightarrow \infty} k(\theta) \quad (2)$$

is uniform in  $\theta \in \Theta$  for every  $s$ , then  $\hat{\theta}_{n_s, n_I}$  is strongly consistent (that is,  $\hat{\theta}_{n_s, n_I}$  converges to  $\theta_0$  almost surely).

Notice that the number of simulations  $n_s$  can remain bounded, it is only  $n_I$  that must tend to infinity for consistence. For practical implementations, it is a crucial point that the  $(U_i^s)$  must be drawn at the beginning of the exploration of the parameter space and kept fixed afterwards while different parameter values are proposed, in order to avoid introducing an extra layer of fluctuation [12, p. 29]. This way, a gradient-based search of the parameter space is possible. At the theoretical level, in the limit  $n_I \rightarrow \infty$ , the estimator is strongly consistent even without using common random numbers.

Under the additional condition that  $\tilde{g}(\mathcal{Y}_i, U_i^s, \theta)$  is twice differentiable with respect to  $\theta$ , asymptotic normality of the estimator also holds and the asymptotic variance can be explicitly given [12, 13].

For the MSM applied to our percolation model with colouring, the data points  $\tilde{Y}_i$  are neither identically distributed (because of boundary effects) nor independent, and Proposition 1 in its current form does not imply the validity of the method. The main theoretical result of this paper is the proof of the strong consistence of a particular MSM estimator for our estimation problem.

The generalised moment function  $K$  we propose contains, in addition to first moments  $Y_i^\ell$ , products  $Y_i^\ell Y_j^\ell$  for  $i \sim j$  because these carry much information about open edges. We note the consequence that it no longer suffices that  $K$  is a function of individual  $Y_i^\ell$  only.

We assume without proof that for this generalised moment function, identifiability (1) holds. For supporting evidence, turn to Section A of the Appendix. This assumption is not true in some extreme cases which we exclude. If  $(\lambda^1, \dots, \lambda^{n_c}) = h \in \{0, 1\}^{n_c}$ , then  $\tilde{Y}$  is almost surely identically  $h$  for any choice of  $\mu$  (and so is  $Y$ ). For an  $h \in \{0, 1\}^{n_c}$ , the outcome  $\tilde{Y}$  is again  $h$  with high probability as  $n_I \rightarrow \infty$ , if  $\mu = 1$ , and  $\lambda^\ell > 0$  if and only if  $h_\ell = 1$ .

The percolation parameter  $\mu$  is allowed to take any value in the subcritical regime  $[0, p_c[$ .  $p_c$  is the *critical probability* of bond percolation. For the triangular lattice, its value is  $p_c = 2 \sin \frac{\pi}{18} \approx 0.3473$ , while for the square lattice, it is  $p_c = 1/2$  [21], [14, Chapter 3].

Section 3 details the steps leading to the SLLN result (2). Due to dependence between the  $Y_i$ , cross-correlations appear in the derivation in addition to variances. Section 4 deals with upper bounding these correlations using percolation theory. Section 5 describes the extension of the SLLN to ULLN.

The observed colouring of the dataset is denoted by  $\mathcal{Y}_i^\ell$  ( $i \in I$ ,  $\ell \in \{1, 2, \dots, n_c\}$ ), whereas in the simulated data it is  $\tilde{Y}_i^{\ell,s}$  ( $s \in \{1, 2, \dots, n_s\}$ ). While it is clear that the simulated data must come from a finite  $I$  (or perhaps from some  $I' : I \subset I' \subset L$ ), we leave flexibility whether the data is of type  $(\tilde{\mathcal{Y}}_i)_{i \in I'}$ , which is the case in our practical application, or of the theoretically appealing type  $(\mathcal{Y}_i)_{i \in I}$ . We let  $(\mathcal{Y}_i)_{i \in I}$  denote both cases, to be interpreted as the context demands. Lastly, we introduce the following averages:

$$\begin{aligned} \bar{\mathcal{Y}}^\ell &:= \frac{1}{n_I} \sum_{i \in I} \mathcal{Y}_i^\ell, & \bar{Y}^{\ell,s} &:= \frac{1}{n_I} \sum_{i \in I} \tilde{Y}_i^{\ell,s}, \\ \bar{\mathcal{Z}}^\ell &:= \frac{1}{n_p} \sum_{(i,j) \in I_2} \mathcal{Y}_i^\ell \mathcal{Y}_j^\ell, & \bar{Z}^{\ell,s} &:= \frac{1}{n_p} \sum_{(i,j) \in I_2} \tilde{Y}_i^{\ell,s} \tilde{Y}_j^{\ell,s}. \end{aligned}$$

Our main theorem is the following.

**Theorem 2** Let  $\Theta$  be a compact subset of  $([0, 1]^{n_c} \setminus \{0, 1\}^{n_c}) \times [0, p_c[$ . (For the triangular lattice,  $p_c = 2 \sin \frac{\pi}{18} \approx 0.3473$ , while in the square lattice case,  $p_c = 1/2$ .) Consider the bond percolation model with colouring and with the true parameter value  $\theta_0 = (\lambda^1, \dots, \lambda^{n_c}, \mu) \in \Theta$ . Let  $\Omega \in \mathbb{R}^{2n_c \times 2n_c}$  be a symmetric, positive definite matrix, which might be a function of the data, and write  $\alpha(\eta) = \eta^T \Omega \eta$  for the resulting quadratic form. Under the assumption of identifiability, when  $n_s$  is fixed and  $n_I$  tends to infinity,

$$\begin{aligned} \hat{\theta}_{n_s, n_I} &:= \arg \min_{\theta \in \Theta} \alpha \left( \begin{pmatrix} \left( \frac{1}{n_I} \sum_{i \in I} \left( \mathcal{Y}_i^\ell - \frac{1}{n_s} \sum_{s=1}^{n_s} \tilde{Y}_i^{\ell,s} \right) \right)_{\ell \in \{1, \dots, n_c\}} \\ \left( \frac{1}{n_p} \sum_{(i,j) \in I_2} \left( \mathcal{Y}_i^\ell \mathcal{Y}_j^\ell - \frac{1}{n_s} \sum_{s=1}^{n_s} \tilde{Y}_i^{\ell,s} \tilde{Y}_j^{\ell,s} \right) \right)_{\ell \in \{1, \dots, n_c\}} \end{pmatrix} \right) \\ &= \arg \min_{\theta \in \Theta} \alpha \left( \begin{pmatrix} \left( \bar{\mathcal{Y}}^\ell - \frac{1}{n_s} \sum_{s=1}^{n_s} \bar{Y}^{\ell,s} \right)_{\ell \in \{1, \dots, n_c\}} \\ \left( \bar{\mathcal{Z}}^\ell - \frac{1}{n_s} \sum_{s=1}^{n_s} \bar{Z}^{\ell,s} \right)_{\ell \in \{1, \dots, n_c\}} \end{pmatrix} \right) \end{aligned}$$

is strongly consistent.

In order to prove the claim, we want to establish that for the arithmetic means generated under general  $\theta$ , the following almost sure convergences hold as  $n_I \rightarrow \infty$ , uniformly in  $\theta \in \Theta$ :

$$\begin{aligned} \frac{1}{n_I} \sum_{i \in I} Y_i^\ell - \frac{1}{n_I} \sum_{i \in I} \mathbb{E}_\theta Y_i^\ell &\rightarrow 0 \\ \text{and } \frac{1}{n_p} \sum_{(i,j) \in I_2} Y_i^\ell Y_j^\ell - \frac{1}{n_p} \sum_{(i,j) \in I_2} \mathbb{E}_\theta [Y_i^\ell Y_j^\ell] &\rightarrow 0, \end{aligned}$$

for  $i \sim j$ . The same proofs apply with  $\tilde{Y}$ , too. This unusual formulation of the SLLN is needed because the random variables  $\tilde{Y}$  are not identically distributed due to boundary effects. These two

SLLNs ultimately ensure that

$$\alpha \left( \begin{array}{c} \left( \bar{\mathcal{Y}}^\ell - \frac{1}{n_s} \sum_{s=1}^{n_s} \bar{Y}^{\ell,s} \right)_{\ell \in \{1, \dots, n_c\}} \\ \left( \bar{\mathcal{Z}}^\ell - \frac{1}{n_s} \sum_{s=1}^{n_s} \bar{Z}^{\ell,s} \right)_{\ell \in \{1, \dots, n_c\}} \end{array} \right) - \alpha \left( \begin{array}{c} \left( \frac{1}{n_I} \sum_{i \in I} \left( E_0 Y_i^\ell - E_\theta \tilde{Y}_i^\ell \right) \right)_{\ell \in \{1, \dots, n_c\}} \\ \left( \frac{1}{n_p} \sum_{(i,j) \in I_2} \left( E_0 [Y_i^\ell Y_j^\ell] - E_\theta [\tilde{Y}_i^\ell \tilde{Y}_j^\ell] \right) \right)_{\ell \in \{1, \dots, n_c\}} \end{array} \right) \xrightarrow{n_I \rightarrow \infty} 0 \quad (3)$$

uniformly with probability 1. The right term is minimal when it is asymptotically zero (in the case when  $E_0$  acts on  $Y_i^\ell$  and  $Y_i^\ell Y_j^\ell$ ; when it acts on  $\tilde{Y}_i^\ell$  and  $\tilde{Y}_i^\ell \tilde{Y}_j^\ell$ , then it is actually zero), and this is achieved in only  $\theta = \theta_0$  under the assumption of identifiability (1). This gives the strong consistence for  $\hat{\theta}_{n_s, n_I}$ .

### 3 Strong law of large numbers with weak dependence

We adapt the proof of Theorem 1 of [8] in this section to suit our purposes. We write out the claims with  $Y$ , but they also hold for  $\tilde{Y}$ .

**Proposition 3** Let  $\theta \in [0, 1]^{n_c} \times [0, p_c[$ , where  $p_c$  is the critical probability of bond percolation. If  $Y$  is generated with parameter value  $\theta$ , then

$$\frac{1}{n_I} \left( \sum_{i \in I} Y_i^\ell - \sum_{i \in I} E_\theta Y_i^\ell \right) \xrightarrow{n_I \rightarrow \infty} 0$$

almost surely. The claim also holds for  $\tilde{Y}$ .

**Proposition 4** Let  $\theta \in [0, 1]^{n_c} \times [0, p_c[$ . If  $Y$  is generated with parameter value  $\theta$ , then

$$\frac{1}{n_p} \left( \sum_{(i,j) \in I_2} Y_i^\ell Y_j^\ell - \sum_{(i,j) \in I_2} E_\theta [Y_i^\ell Y_j^\ell] \right) \xrightarrow{n_I \rightarrow \infty} 0$$

almost surely. The claim also holds for  $\tilde{Y}$ .

**Proof (Proposition 3)** For the ease of notation, let  $Y_i := Y_i^\ell$  for some fixed  $\ell \in \{1, \dots, n_c\}$  ( $i \in I$ ), created by our percolation process with  $\theta = (\lambda^1, \dots, \lambda^{n_c}, \mu)$ . Let  $a > 1$  and define the lacunary sequence  $k_n := [a^n]$ . Let  $S_k := \sum_{i=1}^k Y_i$ .

By the application of Chebyshov's inequality, for every  $\varepsilon > 0$ ,

$$\begin{aligned} \sum_{n=1}^{\infty} P \left( \left| \frac{S_{k_n} - E S_{k_n}}{k_n} \right| > \varepsilon \right) &\leq \sum_{n=1}^{\infty} \frac{\text{Var } S_{k_n}}{\varepsilon^2 k_n^2} \\ &\leq \frac{1}{\varepsilon^2} \sum_{n=1}^{\infty} \frac{1}{k_n^2} \sum_{i=1}^{k_n} \text{Var } Y_i \\ &\quad + \frac{1}{\varepsilon^2} \sum_{n=1}^{\infty} \frac{1}{k_n^2} \sum_{1 \leq i \neq j \leq k_n} (E[Y_i Y_j] - E Y_i E Y_j). \end{aligned} \quad (4)$$

If we can prove that this is finite, then by the Borel–Cantelli lemma, as  $n \rightarrow \infty$ , for every  $\theta \in \Theta$ ,

$$\left| \frac{S_{k_n} - E S_{k_n}}{k_n} \right| \rightarrow 0 \quad \text{a.s.} \quad (5)$$

We first show that

$$\sum_{n=1}^{\infty} \frac{1}{k_n^2} \sum_{i=1}^{k_n} \text{Var } Y_i < \infty$$

by noticing that  $\sup_{i \in I} \text{Var } Y_i \leq 1$  and by the following lemma.

**Lemma 5** If  $1 < a$ , then

$$\sum_{n=1}^{\infty} \frac{1}{k_n} < \infty.$$

**Proof** For  $n \in \mathbb{N}$  sufficiently large,  $a^n/2 \leq k_n$  because  $n \geq \frac{\log 2}{\log a}$  suffices. To see this, consider that  $a^n/2 \leq a^n - 1 < k_n$  is achieved, giving the threshold, if  $2 \leq a^n$ . Let

$$N_0 := \max \left\{ 1, \left\lceil \frac{\log 2}{\log a} \right\rceil \right\}.$$

Consequently, for some constant  $c$ ,

$$\sum_{n=1}^{\infty} \frac{1}{k_n} = \sum_{n=1}^{N_0-1} \frac{1}{k_n} + \sum_{n=N_0}^{\infty} \frac{1}{k_n} \leq c + \sum_{n=N_0}^{\infty} \frac{2}{a^n} = c + \frac{2}{a^{N_0}(1-1/a)} < \infty. \quad \square$$

We prove in Section 4 that

$$\left| \sum_{1 \leq i \neq j \leq k_n} (\mathbb{E}[Y_i Y_j] - \mathbb{E}Y_i \mathbb{E}Y_j) \right| = \mathcal{O}(k_n), \quad (6)$$

so that by applying Lemma 5 once again, we get that (4) is finite, as required.

In the case of a general  $k := n_I$ ,  $k$  is sandwiched between some  $k_n \leq k < k_{n+1}$  and

$$\begin{aligned} \frac{S_k - \mathbb{E}S_k}{k} &\leq \frac{S_{k_{n+1}} - \mathbb{E}S_{k_n}}{k} \\ &\leq \left| \frac{S_{k_{n+1}} - \mathbb{E}S_{k_{n+1}}}{k_{n+1}} \right| \frac{k_{n+1}}{k} + \frac{\mathbb{E}S_{k_{n+1}} - \mathbb{E}S_{k_n}}{k_n}. \end{aligned} \quad (7)$$

Note that even for  $S_{k_{n+1}} - \mathbb{E}S_{k_n} < 0$ , one can change the denominator from  $k$  to  $k_n$  in the second inequality because the right-hand side is nonnegative. Here, for a fixed  $a > 1$ ,

$$\frac{k_{n+1}}{k_n} = \frac{\lceil a^{n+1} \rceil}{\lceil a^n \rceil} \leq \frac{a^{n+1}}{a^n - 1} = a + \frac{a}{a^n - 1}, \quad (8)$$

which in turn is arbitrarily close to  $a$  when  $n$  is sufficiently large. Additionally,

$$\begin{aligned} \frac{\mathbb{E}S_{k_{n+1}} - \mathbb{E}S_{k_n}}{k_n} &\leq \frac{(k_{n+1} - k_n) \sup_{i \in I} \mathbb{E}Y_i}{k_n} \\ &\leq \left( a + \frac{a}{a^n - 1} - 1 \right) \sup_{i \in I} \mathbb{E}Y_i, \end{aligned}$$

and combining this with (5) yields

$$\limsup_{k \rightarrow \infty} \frac{S_k - \mathbb{E}S_k}{k} \leq (a - 1) \sup_{i \in I} \mathbb{E}Y_i \leq a - 1.$$

A similar lower bound can also be attained. Since  $a > 1$  can be chosen arbitrarily, the SLLN for  $Y_i$  (Proposition 3) holds once we prove the estimate (6).  $\square$

**Proof (Proposition 4)** This proof goes entirely analogously to that of Proposition 3. We keep using the notation  $Y_i := Y_i^\ell$  for some fixed  $\ell \in \{1, \dots, n_c\}$  and fixed  $\theta$ , and the lacunary sequence  $k_n = \lceil a^n \rceil$  for  $a > 1$ . Let  $T_k := \sum_{(i,j) \in I_2} Y_i Y_j$  for  $I = I(k)$  composed of the first  $k$  vertices according

to the fixed ordering. This sum has  $n_p(k)$  terms. Then, by the argument of (4), for every  $\varepsilon > 0$ ,

$$\begin{aligned}
\sum_{n=1}^{\infty} \mathbb{P} \left( \left| \frac{T_{k_n} - \mathbb{E}T_{k_n}}{n_p(k_n)} \right| > \varepsilon \right) &\leq \sum_{n=1}^{\infty} \frac{\text{Var } T_{k_n}}{\varepsilon^2 n_p(k_n)^2} \\
&\leq \frac{1}{\varepsilon^2} \sum_{n=1}^{\infty} \frac{1}{n_p(k_n)^2} \sum_{(i_1, i_2) \in I_2(k_n)} \text{Var}[Y_{i_1} Y_{i_2}] \\
&\quad + \frac{1}{\varepsilon^2} \sum_{n=1}^{\infty} \frac{1}{n_p(k_n)^2} \sum_{\substack{(i_1, i_2), (j_1, j_2) \in I_2(k_n) \\ (i_1, i_2) \neq (j_1, j_2)}} \left( \mathbb{E}[Y_{i_1} Y_{i_2} Y_{j_1} Y_{j_2}] \right. \\
&\quad \left. - \mathbb{E}[Y_{i_1} Y_{i_2}] \mathbb{E}[Y_{j_1} Y_{j_2}] \right). \tag{9}
\end{aligned}$$

By

$$\sup_{(i_1, i_2) \in I_2(k_n)} \text{Var}[Y_{i_1} Y_{i_2}] \leq 1$$

and  $|I_2(k_n)| = n_p(k_n) \sim 3n_I = 3k_n$ , Lemma 5 gives

$$\frac{1}{\varepsilon^2} \sum_{n=1}^{\infty} \frac{1}{n_p(k_n)^2} \sum_{(i_1, i_2) \in I_2(k_n)} \text{Var}[Y_{i_1} Y_{i_2}] \leq \frac{1}{\varepsilon^2} \sum_{n=1}^{\infty} \frac{1}{n_p(k_n)} < \infty.$$

In Section 4, it is shown that

$$\left| \sum_{\substack{(i_1, i_2), (j_1, j_2) \in I_2(k_n) \\ (i_1, i_2) \neq (j_1, j_2)}} \left( \mathbb{E}[Y_{i_1} Y_{i_2} Y_{j_1} Y_{j_2}] - \mathbb{E}[Y_{i_1} Y_{i_2}] \mathbb{E}[Y_{j_1} Y_{j_2}] \right) \right| = \mathcal{O}(k_n), \tag{10}$$

and by Lemma 5, we get that the sum (9) is finite. By the Borel–Cantelli lemma,

$$\left| \frac{T_{k_n} - \mathbb{E}T_{k_n}}{n_p(k_n)} \right| \rightarrow 0 \quad \text{a.s.} \tag{11}$$

For a general  $k = n_I$  with  $k_n \leq k < k_{n+1}$ ,

$$\frac{T_k - \mathbb{E}T_k}{n_p(k)} \leq \left| \frac{T_{k_{n+1}} - \mathbb{E}T_{k_{n+1}}}{n_p(k_{n+1})} \right| \frac{n_p(k_{n+1})}{n_p(k)} + \frac{\mathbb{E}T_{k_{n+1}} - \mathbb{E}T_{k_n}}{n_p(k_n)}. \tag{12}$$

For a fixed  $a > 1$ , by using (8) again,

$$\frac{n_p(k_{n+1})}{n_p(k_n)} \sim \frac{3k_{n+1}}{3k_n} \leq a + \frac{a}{a^n - 1},$$

and the right-hand side is arbitrarily close to  $a$  when  $n$  is sufficiently large. Additionally,

$$\frac{\mathbb{E}T_{k_{n+1}} - \mathbb{E}T_{k_n}}{n_p(k_n)} \leq \frac{(n_p(k_{n+1}) - n_p(k_n)) \sup_{(i_1, i_2) \in I_2(k_n)} \mathbb{E}[Y_{i_1} Y_{i_2}]}{n_p(k_n)},$$

hence

$$\limsup_{n \rightarrow \infty} \frac{\mathbb{E}T_{k_{n+1}} - \mathbb{E}T_{k_n}}{n_p(k_n)} \leq (a - 1) \sup_{(i_1, i_2) \in I_2(k_n)} \mathbb{E}[Y_{i_1} Y_{i_2}].$$

Combining this with (11) and (12), we get

$$\limsup_{k \rightarrow \infty} \frac{T_k - \mathbb{E}T_k}{n_p(k)} \leq (a - 1) \sup_{(i_1, i_2) \in I_2(k_n)} \mathbb{E}[Y_{i_1} Y_{i_2}] \leq a - 1.$$

A similar lower bound can also be attained. Since  $a > 1$  can be chosen arbitrarily, the SLLN for  $Y_i Y_j$ ,  $i \sim j$  (Proposition 4) holds once we prove the estimate (10).  $\square$



## 4 Upper bound on correlations

We prove the estimates (6) and (10) in greater generality, for every positive integer  $n$ . Let  $\Theta$  be a compact subset of  $[0, 1] \times [0, p_c]$ , where  $p_c$  is the critical probability of bond percolation.

**Lemma 6** As  $n \rightarrow \infty$ , it holds

$$\sup_{\theta \in \Theta} \left| \sum_{1 \leq i \neq j \leq n} (E[Y_i Y_j] - EY_i EY_j) \right| = \mathcal{O}(n).$$

**Lemma 7** As  $n \rightarrow \infty$ , it holds

$$\sup_{\theta \in \Theta} \left| \sum_{\substack{(i_1, i_2), (j_1, j_2) \in I_2(n) \\ (i_1, i_2) \neq (j_1, j_2)}} \left( E[Y_{i_1} Y_{i_2} Y_{j_1} Y_{j_2}] - E[Y_{i_1} Y_{i_2}] E[Y_{j_1} Y_{j_2}] \right) \right| = \mathcal{O}(n).$$

For background, first we recapitulate from the fundamentals of percolation theory the meaning of increasing events, the FKG inequality, disjoint occurrence, the BK inequality and pivotality [14, Chapter 2]. It is well known that these concepts do not rely on the specific structure of the lattice graph and can be cast more generally in terms of functions of Boolean variables.

In this vein, one can consider a probability space  $(\Gamma, \mathcal{F}, P)$  with sample space  $\Gamma = \{0, 1\}^S$  ( $S$  is finite or at most countably infinite) where the set of events  $\mathcal{F}$  is the  $\sigma$ -algebra generated by the finite-dimensional cylinder sets and the measure is a product measure

$$P = \prod_{s \in S} \nu_s$$

where  $\nu_s$  is specified by some vector  $(p(s))_{s \in S} \in [0, 1]^S$  via

$$\nu_s(\omega(s) = 1) = p(s), \quad \nu_s(\omega(s) = 0) = 1 - p(s)$$

for sample vectors  $(\omega(s))_{s \in S} \in \{0, 1\}^S$  [14, Chapter 2, p. 33].

In our application, we have already fixed a colour  $\ell \in \{1, \dots, n_c\}$  and look at colours independently. We extend the set of vertices  $L$  with an additional vertex that we call  $\infty^\ell$ , or simply  $\infty$  when the colour is fixed and unimportant:  $L^* := L \cup \{\infty\}$ . We also extend the edge set of the triangular lattice  $L_2$  with edges between each vertex and  $\infty$ , and the value assigned to such an edge indicates the presence or absence of seeding. We call these edges *source edges*. For the source edges,  $p(s) = \lambda^\ell$ , and for the edges of the lattice which represent contamination,  $p(s) = \mu$ . The interpretation is that  $Y_i^\ell = 1$  if and only if  $i \leftrightarrow^* \infty^\ell$ , where the asterisk refers to connection in the extended graph.

An event  $A \in \mathcal{F}$  of the  $\sigma$ -algebra is called *increasing*, if whenever  $\omega \leq \omega'$ ,  $\omega \in A$  implies  $\omega' \in A$ .

**Theorem 8 (FKG inequality [9],[14, pp. 34–36])** If  $A$  and  $B$  are increasing events in  $\mathcal{F}$ , then  $P(A \cap B) \geq P(A)P(B)$ .

Let  $e_1, e_2, \dots, e_N$  be  $N$  distinct edges of the graph, and  $A, B \in \mathcal{F}$  two increasing events which depend on the vector of the states of these  $N$  edges  $\omega = (\omega(e_1), \dots, \omega(e_N))$  only. Such vectors  $\omega$  are characterised uniquely by the set of edges with value 1:  $J(\omega) = \{e_i \mid i \in \{1, \dots, N\}, \omega(e_i) = 1\}$ .

For the increasing events  $A, B$ , the event  $A \circ B$  (we say *A and B occur disjointly*) is the set of all  $\omega \in \Gamma$  for which there exists an  $H \subseteq J(\omega)$  such that  $\omega'$  determined by  $J(\omega') = H$  belongs to  $A$ , and  $\omega''$  determined by  $J(\omega'') = J(\omega) \setminus H$  belongs to  $B$ . In words,  $A \circ B$  is the set of assignments of 0 and 1 to the edges for which there exist two disjoint sets of edges assigned the value 1 (*open edges*) such that the first such set ensures the occurrence of event  $A$  and the second set ensures the occurrence of  $B$ . It is easy to verify that  $A \circ B$  is also increasing and  $A \circ B \subseteq A \cap B$ .

The classical example for disjoint occurrence is when  $A$  is the event that there is an open path joining  $i_1$  to  $j_1$  within the finite subgraph given by  $\{e_1, \dots, e_N\}$  and  $B$  is the event that there is an open path between  $i_2$  and  $j_2$  within the same finite subgraph. Then  $A \circ B$  is the event that there exist two edge-disjoint paths, the first between  $i_1$  and  $j_1$  and another one joining  $i_2$  to  $j_2$ .

**Theorem 9 (BK inequality [4],[14, pp. 37–41])** If  $A$  and  $B$  are increasing events in  $\mathcal{F}$ , then  $P(A \circ B) \leq P(A)P(B)$ .

The validity of the inequality extends to the existence of arbitrarily long (but finite length) edge-disjoint open paths, which is what we need it for, by taking a sequence of growing, nested subsets of  $L$  [14, p. 38].

The notion of pivotality is not used until Section 5. For any event  $A$  an edge  $e$  is pivotal if its open or closed state is crucial to whether  $A$  occurs or not. In more detail, the edge  $e$  is pivotal for the pair  $(A, \omega)$ , if for the indicator function of  $A$ ,  $\chi_A(\omega) \neq \chi_A(\omega')$ , where the configuration  $\omega' \in \{0, 1\}^S$  is defined by  $\omega'(e) = 1 - \omega(e)$ , and  $\omega'(f) = \omega(f)$  for every edge  $f \neq e$ . The event that  $e$  is pivotal for  $A$  is the set of  $\omega$  for which  $e$  is pivotal for  $(A, \omega)$ .

**Proof (Lemma 6)** In the extended lattice graph that has source edges with weight zero or one at every vertex for seeding, the event  $\{Y_i = 1\}$  for  $i \in I$  is increasing because it is increasing in both seeding (source edges) and contamination edges. For any  $i, j \in I$ ,

$$E[Y_i Y_j] - EY_i EY_j = P(Y_i Y_j = 1) - P(Y_i = 1)P(Y_j = 1) \geq 0$$

by the FKG inequality. Hence, for every  $\theta \in \Theta$ ,

$$\sum_{1 \leq i \neq j \leq n} (E[Y_i Y_j] - EY_i EY_j) \geq 0.$$

For the upper bound, consider that

$$\begin{aligned} P(Y_i Y_j = 1) - P(Y_i = 1)P(Y_j = 1) &= P(\{Y_i = 1\} \circ \{Y_j = 1\}) - P(Y_i = 1)P(Y_j = 1) \\ &\quad + P(\{Y_i Y_j = 1\} \setminus \{Y_i = 1\} \circ \{Y_j = 1\}) \\ &\leq P(\{Y_i Y_j = 1\} \setminus \{Y_i = 1\} \circ \{Y_j = 1\}) \end{aligned} \quad (13)$$

by the BK inequality. Cooccurrence of  $\{Y_i = 1\}$  and  $\{Y_j = 1\}$  which is not disjoint is one where  $i$  and  $j$  are in the same component in the edge set on the non-extended lattice:

$$\{Y_i Y_j = 1\} \setminus \{Y_i = 1\} \circ \{Y_j = 1\} \subseteq \{i \leftrightarrow j\}.$$

We show that

$$\sum_{1 \leq i \neq j \leq n} P(i \leftrightarrow j) = \mathcal{O}(n) \quad (14)$$

for  $\mu < p_c$ , and uniformly so for  $\mu \in [0, p_c - \varepsilon]$  for every  $\varepsilon > 0$ . This follows from the exponential decay of the cluster size distribution and it will complete the proof of Lemma 6.

Let  $C(i)$  denote the set of vertices in the component of  $i \in L$  according to the non-extended edge set of  $L$ . Then

$$\sum_{1 \leq i \neq j \leq n} P(i \leftrightarrow j) = \sum_{1 \leq i \leq n} \sum_{\substack{1 \leq j \leq n \\ j \neq i}} E\chi_{\{i \leftrightarrow j\}} = \sum_{1 \leq i \leq n} E[|C(i)| - 1].$$

**Theorem 10 (Exponential decay of the cluster size distribution [1], [14, Chapter 6])** For  $\mu \in ]0, p_c[$ , there exists  $g(\mu) > 0$  such that for all  $k \geq 1$  and  $i \in I$ , for the bond percolation with parameter  $\mu$ , it holds that  $P(|C(i)| \geq k) \leq e^{-kg(\mu)}$ .

Take  $\mu^* = p_c - \varepsilon$ . As  $P(|C(i)| \geq k)$  is nondecreasing in  $\mu$ , we get a uniform bound in  $\theta \in \Theta$  if the bound is valid for  $\mu^*$ :

$$\begin{aligned} \sum_{i=1}^n \mathbb{E}[|C(i)| - 1] &= \sum_{i=1}^n \left( \left( \sum_{k=1}^{\infty} P(|C(i)| \geq k) \right) - 1 \right) \\ &\leq \sum_{i=1}^n \sum_{k=1}^{\infty} e^{-kg(\mu^*)} = n \frac{1}{e^{g(\mu^*)} - 1}. \end{aligned} \quad (15)$$

This proves Lemma 6, which in turn completes the proof of Proposition 3.  $\square$

To go from the case of  $Y$  to  $\tilde{Y}$ , first we couple the realisations of  $(Y_i)_{i \in L}$  and  $(\tilde{Y}_i)_{i \in I}$  with varying lattices  $I$  (and later with varying parameter vectors) by defining them via shared random variables  $(U_i^\ell)_{i \in L, \ell \in \{1, 2, \dots, n_c\}}$  and  $(V_{ij})_{(i,j) \in L_2}$  that are independent and all uniformly distributed on  $[0, 1]$ . For  $\theta = (\lambda^1, \dots, \lambda^{n_c}, \mu) \in \Theta$ ,  $i \in L$ ,  $(i, j) \in L_2$  and  $\ell \in \{1, 2, \dots, n_c\}$ , the seeding is defined by  $X_i^\ell := \chi_{\{U_i^\ell < \lambda^\ell\}}$ , and edges are open according to  $\xi_{ij} := \chi_{\{V_{ij} < \mu\}}$ .

Let us drop the superscript  $\ell$  again. Notice that any  $\tilde{Y}_i$  can increase when  $I$  is increased. In the proof of Proposition 3, the only occasion where  $Y_i$  from different  $I$  are compared is inequality (7). We mark the lattice size as a variable in the superscript of  $\tilde{Y}_i$ . Observe that  $\tilde{Y}_i^{k_n} \leq \tilde{Y}_i^k \leq \tilde{Y}_i^{k_{n+1}} \leq Y_i$  for  $k_n \leq k < k_{n+1}$  and  $i \in I(k_n)$ . With  $\tilde{S}_n^k := \sum_{i=1}^n \tilde{Y}_i^k$ , noting  $\tilde{S}_{k_n}^{k_n} \leq \tilde{S}_k^k \leq \tilde{S}_{k_{n+1}}^{k_{n+1}}$ ,

$$\begin{aligned} \frac{\tilde{S}_k^k - \mathbb{E}\tilde{S}_k^k}{k} &\leq \frac{\tilde{S}_{k_{n+1}}^{k_{n+1}} - \mathbb{E}\tilde{S}_{k_n}^{k_n}}{k} \\ &\leq \left| \frac{\tilde{S}_{k_{n+1}}^{k_{n+1}} - \mathbb{E}\tilde{S}_{k_{n+1}}^{k_{n+1}}}{k_{n+1}} \right| \frac{k_{n+1}}{k_n} + \frac{\mathbb{E}\tilde{S}_{k_{n+1}}^{k_{n+1}} - \mathbb{E}\tilde{S}_{k_n}^{k_n}}{k_n}, \end{aligned}$$

where, similarly to inequality (7), the second inequality holds for different reasons when  $\tilde{S}_{k_{n+1}}^{k_{n+1}} - \mathbb{E}\tilde{S}_{k_n}^{k_n}$  is negative and when not. The first term does not require special treatment. The second term is

$$\frac{\mathbb{E}\tilde{S}_{k_{n+1}}^{k_{n+1}} - \mathbb{E}\tilde{S}_{k_n}^{k_n}}{k_n} = \frac{1}{k_n} \sum_{i=k_n+1}^{k_{n+1}} \mathbb{E}\tilde{Y}_i^{k_{n+1}} + \frac{1}{k_n} \sum_{i=1}^{k_n} (\mathbb{E}\tilde{Y}_i^{k_{n+1}} - \mathbb{E}\tilde{Y}_i^{k_n}). \quad (16)$$

Here

$$\frac{1}{k_n} \sum_{i=k_n+1}^{k_{n+1}} \mathbb{E}\tilde{Y}_i^{k_{n+1}} \leq \frac{k_{n+1} - k_n}{k_n} \sup_{i \in I(k_{n+1})} \mathbb{E}\tilde{Y}_i^{k_{n+1}}$$

is dealt with as in the original proof of Proposition 3. For the other term of (16),

$$\frac{1}{k_n} \sum_{i=1}^{k_n} (\mathbb{E}\tilde{Y}_i^{k_{n+1}} - \mathbb{E}\tilde{Y}_i^{k_n}) \leq \frac{1}{k_n} \sum_{i=1}^{k_n} (\mathbb{E}Y_i - \mathbb{E}\tilde{Y}_i^{k_n}).$$

According to the next proposition, this vanishes in the limit, leaving us with

$$\limsup_{k \rightarrow \infty} \frac{\tilde{S}_k^k - \mathbb{E}\tilde{S}_k^k}{k} \leq a - 1,$$

as required.

**Proposition 11** For a compact subset  $\Theta \subset [0, 1] \times [0, p_c]$ ,

$$\sup_{\theta \in \Theta} \frac{1}{n_I} \left| \sum_{i \in I} \mathbb{E}_\theta \tilde{Y}_i - \sum_{i \in I} \mathbb{E}_\theta Y_i \right| \xrightarrow{n_I \rightarrow \infty} 0.$$

**Proof** As  $Y_i \geq \tilde{Y}_i$  almost surely,

$$\begin{aligned} \mathbb{E}[Y_i - \tilde{Y}_i] &= \mathbb{P}(Y_i = 1, \tilde{Y}_i = 0) \\ &\leq \mathbb{P}(Y_i = 1 \text{ and } \exists j \in L \setminus I : X_j = 1, i \leftrightarrow j) \\ &\leq \mathbb{P}(i \leftrightarrow \Delta I), \end{aligned}$$

which expresses that  $Y_i$  and  $\tilde{Y}_i$  can differ only if  $i \in I$  is connected to the exterior vertex boundary of  $I$ . Further,

$$\sum_{i \in I} \mathbb{P}(i \leftrightarrow \Delta I) = \mathbb{E} \left[ \sum_{i \in I} \chi_{\{i \leftrightarrow \Delta I\}} \right].$$

But this is the expected size of the open component that is grown from all vertices of  $\Delta I$  towards the inside of  $I$ . It is upper bounded by  $|\Delta I| \times \mathbb{E}[|C(0)|]$ . On the compact  $\Theta$ , the mean size of the open component of any vertex has a universal finite upper bound by (15). Hence,

$$\sup_{\theta \in \Theta} \frac{1}{n_I} \left| \sum_{i \in I} \mathbb{E}_\theta [\tilde{Y}_i - Y_i] \right| \leq \frac{|\Delta I|}{n_I} \mathbb{E}[|C(0)|] \rightarrow 0$$

as  $n_I \rightarrow \infty$ , due to our assumption  $|\Delta I|/|I| \rightarrow 0$  about the nested sequence of  $I$ .  $\square$

As in the case of  $Y$ , a lower bound for  $(\tilde{S}_k^k - \mathbb{E}\tilde{S}_k^k)/k$  does not pose any additional difficulty. In the proof of Lemma 6, the covariances cannot increase when we constrain the set of edges to those among the first  $n$  vertices. Concretely,  $\mathbb{E}[|C(i)|]$  cannot increase. Therefore Proposition 3 stays true for  $\tilde{Y}$ .

**Proof (Lemma 7)** The proof follows closely that of Lemma 6. For any two pairs  $(i_1, i_2), (j_1, j_2) \in I_2$ ,

$$\mathbb{E}[Y_{i_1} Y_{i_2} Y_{j_1} Y_{j_2}] - \mathbb{E}[Y_{i_1} Y_{i_2}] \mathbb{E}[Y_{j_1} Y_{j_2}] \geq 0$$

due to the FKG inequality applied to  $\{Y_{i_1} Y_{i_2} = 1\}$  and  $\{Y_{j_1} Y_{j_2} = 1\}$ . Therefore, for any  $\theta \in \Theta$ ,

$$\sum_{\substack{(i_1, i_2), (j_1, j_2) \in I_2(n) \\ (i_1, i_2) \neq (j_1, j_2)}} \left( \mathbb{E}[Y_{i_1} Y_{i_2} Y_{j_1} Y_{j_2}] - \mathbb{E}[Y_{i_1} Y_{i_2}] \mathbb{E}[Y_{j_1} Y_{j_2}] \right) \geq 0.$$

The first step towards the upper bound, similarly to (13), uses the BK inequality:

$$\begin{aligned} \mathbb{P}(Y_{i_1} Y_{i_2} Y_{j_1} Y_{j_2} = 1) - \mathbb{P}(Y_{i_1} Y_{i_2} = 1) \mathbb{P}(Y_{j_1} Y_{j_2} = 1) \\ \leq \mathbb{P}(\{Y_{i_1} Y_{i_2} Y_{j_1} Y_{j_2} = 1\} \setminus \{Y_{i_1} Y_{i_2} = 1\} \circ \{Y_{j_1} Y_{j_2} = 1\}). \end{aligned}$$

Cooccurrence which is not disjoint is one where at least one of  $i_1$  and  $i_2$  is connected to at least one of  $j_1$  and  $j_2$  in the non-extended edge set, or in symbols,

$$\{Y_{i_1} Y_{i_2} Y_{j_1} Y_{j_2} = 1\} \setminus \{Y_{i_1} Y_{i_2} = 1\} \circ \{Y_{j_1} Y_{j_2} = 1\} \subseteq \{i_1 \leftrightarrow \{j_1, j_2\}\} \cup \{i_2 \leftrightarrow \{j_1, j_2\}\},$$

where  $\leftrightarrow$  is meant to be reflective so that not disjoint cooccurrence might involve e.g. that  $i_1 = j_1$ . So for every fixed  $\theta \in \Theta$ ,

$$\begin{aligned} \sum_{\substack{(i_1, i_2), (j_1, j_2) \in I_2(n) \\ (i_1, i_2) \neq (j_1, j_2)}} \left( \mathbb{E}[Y_{i_1} Y_{i_2} Y_{j_1} Y_{j_2}] - \mathbb{E}[Y_{i_1} Y_{i_2}] \mathbb{E}[Y_{j_1} Y_{j_2}] \right) \\ \leq \sum_{\substack{(i_1, i_2), (j_1, j_2) \in I_2(n) \\ (i_1, i_2) \neq (j_1, j_2)}} \left( \mathbb{P}(i_1 \leftrightarrow \{j_1, j_2\}) + \mathbb{P}(i_2 \leftrightarrow \{j_1, j_2\}) \right). \end{aligned}$$

It suffices to treat the two terms individually, and one of them gives

$$\begin{aligned}
\sum_{\substack{(i_1, i_2), (j_1, j_2) \in I_2(n) \\ (i_1, i_2) \neq (j_1, j_2)}} P(i_1 \leftrightarrow \{j_1, j_2\}) &\leq \sum_{i_1 \in I(n)} \sum_{i_2 \sim i_1} \sum_{(j_1, j_2) \in I_2(n) \setminus \{(i_1, i_2)\}} \left( P(i_1 \leftrightarrow j_1) + P(i_1 \leftrightarrow j_2) \right) \\
&\leq \sum_{i_1 \in I(n)} \sum_{i_2 \sim i_1} \left( \sum_{\substack{(j_1, j_2) \in I_2(n) \setminus \{(i_1, i_2)\} \\ j_1 = i_1}} 2 + \sum_{\substack{(j_1, j_2) \in I_2(n) \setminus \{(i_1, i_2)\} \\ j_2 = i_1}} 2 \right. \\
&\quad \left. + \sum_{(j_1, j_2) \in (I(n) \setminus \{i_1\})_2} \left( P(i_1 \leftrightarrow j_1) + P(i_1 \leftrightarrow j_2) \right) \right) \\
&\leq 6 \times 2 \sum_{i_1 \in I(n)} \left( 5 \times 2 + 6 \sum_{j_1 \in I(n) \setminus \{i_1\}} P(i_1 \leftrightarrow j_1) \right) \\
&\leq 12n \left( 10 + 6 \frac{1}{e^{g(\mu^*)} - 1} \right) = \mathcal{O}(n),
\end{aligned}$$

where in the second inequality, we separate between cases when  $i_1 = j_1$  or  $i_1 = j_2$  and when not, and notice that when they are not equal, then all  $(j_1, j_2)$  pairs are disjoint from  $i_1$ . In the third inequality, we replace the sum for  $i_2 \sim i_1$  by a factor of 6 (for the triangular lattice), and instead of  $(j_1, j_2)$ , we sweep for  $j_1$ , and then for its at most 6 neighbours  $j_2$  separately. Thereby we reduced the problem to the previous case and the fourth inequality follows by (15). This completes the proofs of Lemma 7 and Proposition 4.  $\square$

The proof of Proposition 4 can be adapted to  $\tilde{Y}$ . For example, the following variant of Proposition 11 also holds:

$$\sup_{\theta \in \Theta} \frac{1}{n_p} \left| \sum_{(i, j) \in I_2} E_{\theta} [\tilde{Y}_i \tilde{Y}_j] - \sum_{(i, j) \in I_2} E_{\theta} [Y_i Y_j] \right| \xrightarrow{n_I \rightarrow \infty} 0. \quad (17)$$

## 5 Uniform law of large numbers (ULLN) for our process

In the interests of conciseness, we continue assuming that there is only one colour:  $n_c = 1$ . This leads to no loss of generality. We prove that the SLLNs, Propositions 3 and 4, hold uniformly over the compact parameter set  $\Theta \subset [0, 1] \times [0, p_c]$ . Similarly to the preceding, we write everything out for  $Y$ , but the result is also valid for  $\tilde{Y}$ .

To prove the uniform version of Proposition 3, we check that the conditions of the following theorem hold, where we adapted [19, p. 8, 2 Theorem] or [22, p. 25, Lemma 3.1] to our setting. For the rewriting of the theorem, we exploited that for a sample  $((U_i)_{i \in L}, (V_{ij})_{(i, j) \in L_2})$  of the seeds and edges, any  $Y_i$  is nondecreasing in both  $\lambda$  and  $\mu$ .

**Theorem 12 (cf. [19, p. 8, 2 Theorem], [22, p. 25, Lemma 3.1])** Suppose that for every  $\varepsilon > 0$  there exists a finite set of pairs of parameter vectors

$$\mathcal{P} = \left\{ (\theta_r^L, \theta_r^U) \in ([0, 1] \times [0, p_c])^2 \mid r \in \{1, \dots, N(\varepsilon)\} \right\}$$

such that

1. for every  $r \in \{1, \dots, N(\varepsilon)\}$ , the SLLN holds for  $\theta_r^L$  and  $\theta_r^U$ ; that is, if  $Y$  is generated with parameter value  $\theta_r^L$ , then

$$\frac{1}{n_I} \left( \sum_{i \in I} Y_i - \sum_{i \in I} E_{\theta_r^L} Y_i \right) \xrightarrow{n_I \rightarrow \infty} 0$$

almost surely, and similarly for  $\theta_r^U$ ;

2. for every  $\theta \in \Theta$ , there is an  $r \in \{1, \dots, N(\varepsilon)\}$  such that  $\theta_r^L \leq \theta \leq \theta_r^U$  coordinatewise;
3. for every  $r \in \{1, \dots, N(\varepsilon)\}$  and  $i \in I$ ,  $E_{\theta_r^U} Y_i - E_{\theta_r^L} Y_i \leq \varepsilon$ .

Then the ULLN holds, that is,

$$\sup_{\theta \in \Theta} \frac{1}{n_I} \left| \sum_{i \in I} Y_i - \sum_{i \in I} E_{\theta} Y_i \right| \xrightarrow{n_I \rightarrow \infty} 0$$

almost surely, where  $Y$  is generated with parameter value  $\theta$ .

We construct  $\mathcal{P}$  such that the rectangles  $R_r$  spanned by  $\theta_r^L = (\lambda_r^L, \mu_r^L)$  and  $\theta_r^U = (\lambda_r^U, \mu_r^U)$ , that is, the closed rectangles  $\{(\lambda_r^L, \mu_r^L), (\lambda_r^U, \mu_r^L), (\lambda_r^U, \mu_r^U), (\lambda_r^L, \mu_r^U)\}$ , cover  $\Theta$ . By this construction, Condition 2 holds. No matter how we choose finitely many pairs  $(\theta_r^L, \theta_r^U)$ , Condition 1 holds for each by Proposition 3. We achieve Condition 3 by proving Lipschitz continuity of the expectation  $E_{\theta} Y_i$  in  $\theta$ .

**Lemma 13** For any  $\mu^* \in ]0, p_c[$ , The expectation  $E_{\theta} Y_i$  is Lipschitz continuous in  $\theta$  over the set  $[0, 1] \times [0, \mu^*]$  with some Lipschitz constant  $L_0$ , which is universal for  $i \in L$ .

Lemma 13 guarantees uniform continuity in  $\theta$ . Instead of  $\delta > 0$  that corresponds to the  $\varepsilon$  required by Theorem 12, we take a  $\delta' \in ]0, \delta[$ . For instance,  $\delta' = \varepsilon/(2L_0)$  is suitable. We cover  $\Theta$  with open rectangles  $R_u$  of the above form with diameter  $|\theta_r^U - \theta_r^L| = \delta'$ . We intersect each  $R_u$  with  $[0, 1] \times [0, p_c[$  to avoid overhangs; they remain relatively open. Because of compactness, there is a finite subcover of  $\Theta$  with such potentially trimmed open rectangles. We define  $\mathcal{P}$  via the vertices of these finitely many rectangles.  $R_r$  are now closed rectangles with diameter  $\delta'$ . They are contained in rectangles with diameter  $\delta$ , and Condition 3 is satisfied. In conclusion, a proof of Lemma 13 proves the SLLN for  $Y$ .

**Proof (Lemma 13)** Consider  $\theta, \theta' \in [0, 1] \times [0, \mu^*]$ ,  $\theta = (\lambda, \mu)$  and  $\theta' = (\lambda', \mu')$ . We can assume that  $\theta \leq \theta'$  coordinatewise. If this were not the case, we would prove the inequality for  $\theta^L := (\lambda \wedge \lambda', \mu \wedge \mu')$  and  $\theta^U := (\lambda \vee \lambda', \mu \vee \mu')$ . This suffices since  $|\theta - \theta'| = |\theta^L - \theta^U|$ , and both  $E_{\theta} Y_i$  and  $E_{\theta'} Y_i$  are contained in  $[E_{\theta^L} Y_i, E_{\theta^U} Y_i]$  due to monotonicity.

We identify the vertices of  $L$  with the source edges, and fix an ordering of all source and contamination edges:  $L \cup L_2 = \{e_0, e_1, e_2, \dots\}$ . Let  $\vartheta : \mathbb{N} \rightarrow [0, 1]$  be such that  $\vartheta_k = \lambda$  if  $e_k$  is a source edge, and  $\vartheta_k = \mu$  if  $e_k$  is a contamination edge. Define  $\vartheta'$  analogously with  $\lambda', \mu'$  in place of  $\lambda, \mu$ , respectively. Finally, let  $\theta^k, \theta'^k : \mathbb{N} \rightarrow [0, 1]$  be defined by

$$(\theta^k)_j := \begin{cases} \vartheta'_j & \text{if } j < k, \\ \vartheta_j & \text{if } j \geq k, \end{cases} \quad (\theta'^k)_j := \begin{cases} \vartheta'_j & \text{if } j \leq k, \\ \vartheta_j & \text{if } j > k, \end{cases}$$

for  $k, j \in \mathbb{N}$ . Let  $\omega_{\theta}$  be the configuration that is specified by  $((U_i)_{i \in L}, (V_{ij})_{(i,j) \in L_2})$  and parameter  $\theta$  via  $(X_i)_{i \in L}$  and  $(\xi_{ij})_{(i,j) \in L_2}$ . Then

$$\begin{aligned} E_{\theta'} Y_i - E_{\theta} Y_i &= P_{\theta'}(Y_i = 1) - P_{\theta}(Y_i = 1) \\ &= \sum_{k=0}^{\infty} P(Y_i(\omega_{\theta'^k}) = 1 \text{ and } Y_i(\omega_{\theta^k}) = 0) \\ &= \sum_{k=0}^{\infty} (\vartheta'_k - \vartheta_k) P_{\theta^k}(e_k \text{ is pivotal for } Y_i = 1), \end{aligned}$$

where the second equality is just the law of total probability when we know that  $\{Y_i = 1\}$  is an increasing event, and the third equality is elaborated in [14, pp. 41–43] as such a step is used in the

proof of Russo's formula. Note that the concerns in that derivation related to an infinite edge set do not apply here because we have always got only one edge  $e_k$  whose parameter differs between  $\theta^k$  and  $\theta'^k$ . (The price we pay is that each pivotality is with a different parameter vector  $\theta^k$ .) If  $e_k$  is a source edge, then  $\vartheta'_k - \vartheta_k = \lambda' - \lambda$ , and if  $e_k$  is a contamination edge, then  $\vartheta'_k - \vartheta_k = \mu' - \mu$ . Further,

$$P_{\theta^k}(e_k \text{ is pivotal for } Y_i = 1) \leq \begin{cases} 1, & \text{if } e_k \text{ is the source edge of vertex } i, \\ P_{\theta^k}(j \leftrightarrow i), & \text{if } e_k \text{ is the source edge of vertex } j \neq i, \\ 1, & \text{if } e_k \text{ is an edge incident with } i, \\ P_{\theta^k}(i_1 \leftrightarrow i) + P_{\theta^k}(i_2 \leftrightarrow i), & \text{if } e_k \text{ is the edge } (i_1, i_2), i \neq i_1, i_2. \end{cases}$$

Then

$$\begin{aligned} \sum_{e_k \text{ source edge}} (\vartheta'_k - \vartheta_k) P_{\theta^k}(e_k \text{ is pivotal for } Y_i = 1) \\ \leq (\lambda' - \lambda) \left( 1 + \sum_{j \in L \setminus \{i\}} P_{\theta^k}(j \leftrightarrow i) \right) \\ \leq (\lambda' - \lambda) \left( 1 + \frac{1}{e^{g(\mu^*)} - 1} \right) \end{aligned}$$

by (15). Using tricks from the proof of Lemma 7,

$$\begin{aligned} \sum_{e_k \text{ contamination edge}} (\vartheta'_k - \vartheta_k) P_{\theta^k}(e_k \text{ is pivotal for } Y_i = 1) \\ \leq (\mu' - \mu) \left( 6 + \sum_{(i_1, i_2) \in (L \setminus \{i\})_2} (P_{\theta^k}(i_1 \leftrightarrow i) + P_{\theta^k}(i_2 \leftrightarrow i)) \right) \\ \leq (\mu' - \mu) \left( 6 + 2 \times 6 \sum_{j \in L \setminus \{i\}} P_{\theta^k}(j \leftrightarrow i) \right) \\ \leq (\mu' - \mu) \left( 6 + \frac{12}{e^{g(\mu^*)} - 1} \right). \end{aligned}$$

Consequently,

$$\begin{aligned} E_{\theta'} Y_i - E_{\theta} Y_i &\leq (\lambda' - \lambda) \left( 1 + \frac{1}{e^{g(\mu^*)} - 1} \right) + (\mu' - \mu) \left( 6 + \frac{12}{e^{g(\mu^*)} - 1} \right) \\ &\leq (\lambda' - \lambda + \mu' - \mu) \left( 6 + \frac{12}{e^{g(\mu^*)} - 1} \right) \leq L_0 |\theta' - \theta| \end{aligned}$$

for some  $L_0 > 0$  because in finite dimensions, all norms are equivalent.  $\square$

Lemma 13 for  $E_{\theta}[Y_i Y_j]$  ( $i \sim j$ ) can be shown by a now straightforward adjustment of the original proof. This then implies that the following modification of Theorem 12 holds.

**Theorem 14** Suppose that the conditions of Theorem 12 hold with the following updates to points 1 and 3:

1' for every  $r \in \{1, \dots, N(\varepsilon)\}$ , if  $Y$  is generated with parameter value  $\theta_r^L$ , then

$$\frac{1}{n_p} \left( \sum_{(i,j) \in I_2} Y_i Y_j - \sum_{(i,j) \in I_2} E_{\theta_r^L}[Y_i Y_j] \right) \xrightarrow{n_I \rightarrow \infty} 0$$

almost surely, and similarly for  $\theta_r^U$ ;

3' for every  $r \in \{1, \dots, N(\varepsilon)\}$  and  $(i, j) \in I_2$ ,  $\mathbb{E}_{\theta_r^U}[Y_i Y_j] - \mathbb{E}_{\theta_r^L}[Y_i Y_j] \leq \varepsilon$ .

Then the ULLN holds, that is

$$\sup_{\theta \in \Theta} \frac{1}{n_p} \left| \sum_{(i,j) \in I_2} Y_i Y_j - \sum_{(i,j) \in I_2} \mathbb{E}_\theta[Y_i Y_j] \right| \xrightarrow{n_I \rightarrow \infty} 0$$

almost surely, where  $Y$  is generated with parameter value  $\theta$ .

The derivations and results of this section hold with  $\tilde{Y}$ , too. The assumption of identifiability, Proposition 11 and (17) together guarantee that the second term of (3) converges to zero almost surely if and only if  $\theta = \theta_0$ . Additionally, the arguments of the first term of (3) converge uniformly to those of the second term almost surely, due to the conclusions of Theorems 12 and 14. This proves our main theorem, Theorem 2.

We followed the philosophy that the dataset  $\mathcal{Y}$  comes from the infinite lattice  $L$  although only a finite subset is observed. This is an idealised view that assumes the existence of a process on the infinite lattice. Otherwise, when the dataset is of type  $\tilde{\mathcal{Y}}$ , the derivation is simpler because Proposition 11 and (17) are not needed.

## 6 Computer testing of the proposed method

### 6.1 Implementation

We implemented the proposed MSM parameter estimator in the MATLAB software (The MathWorks, Inc.), and we report our findings in this section. See also [3] for an early version with  $n_c = 3$  colours. For the the objective function

$$\alpha \left( \begin{pmatrix} \bar{\mathcal{Y}}^\ell - \frac{1}{n_s} \sum_{s=1}^{n_s} \bar{Y}^{\ell,s} \\ \bar{\mathcal{Z}}^\ell - \frac{1}{n_s} \sum_{s=1}^{n_s} \bar{Z}^{\ell,s} \end{pmatrix}_{\ell \in \{1, \dots, n_c\}} \right), \quad (18)$$

we chose the quadratic form  $\alpha(\eta) = \eta^T \Omega \eta$  the following way:

$$\Omega = \text{diag} \left( (\bar{\mathcal{Y}}^1)^{-2}, \dots, (\bar{\mathcal{Y}}^{n_c})^{-2}, (\bar{\mathcal{Z}}^1)^{-2}, \dots, (\bar{\mathcal{Z}}^{n_c})^{-2} \right). \quad (19)$$

In the unlikely case that a  $\bar{\mathcal{Y}}^\ell$  or a  $\bar{\mathcal{Z}}^\ell$  is zero, the corresponding diagonal element of  $\Omega$  is set to 1. Through this normalisation, we expect each coordinate to contribute roughly equally to the sum.

Common random numbers are used during the exploration of the parameter space. This removes an element of fluctuation as different  $\theta = (\lambda^1, \dots, \lambda^{n_c}, \mu) \in \Theta$  are tested. We propose two alternative methods for sampling synthetic datasets. Method 1 is the canonical approach. We draw and fix independent random variables from the uniform distribution on  $[0, 1]$ :  $(U_i^{\ell,s})$  for  $\ell \in \{1, \dots, n_c\}$ ,  $s \in \{1, \dots, n_s\}$ ,  $i \in I$ , and  $(V_{ij}^s)$  for  $s \in \{1, \dots, n_s\}$ ,  $(i, j) \in I_2$ . Thereafter, for each parameter vector, seeding and the open or closed state of edges are defined by

$$X_i^{\ell,s} := \begin{cases} 1 & \text{if } U_i^{\ell,s} < \lambda^\ell, \\ 0 & \text{otherwise,} \end{cases} \quad \text{for } \ell \in \{1, \dots, n_c\}, s \in \{1, \dots, n_s\}, i \in I;$$

$$\xi_{ij}^s := \begin{cases} 1 & \text{if } V_{ij}^s < \mu, \\ 0 & \text{otherwise,} \end{cases} \quad \text{for } s \in \{1, \dots, n_s\}, (i, j) \in I_2.$$

This method gives a binomially distributed number of open edges and, similarly, seeded vertices for each colour  $\ell$ .



We anticipate that it is beneficial for the parameter estimation to remove the randomness in the numbers of seeds and open edges, and to make exactly as many edges open as their expected number,  $\zeta(\mu n_p)$ , where  $\zeta$  is the rounding to the nearest integer with some tie-breaking rule. The same is stipulated for seeds:  $\zeta(\lambda^\ell n_I)$  random vertices shall be seeded with colour  $\ell$ . This is what Method 2 does. We see this as a variance-reduction trick that achieves lower variance by introducing dependencies between random draws: for example, by knowing the state of all edges but one, we can infer the state of the remaining edge.

Let  $S_n$  denote the set of permutations of  $\{1, \dots, n\}$ . In Method 2, one draws permutations  $(\sigma^{\ell,s})$  from  $S_{n_I}$  independently, uniformly at random for  $\ell \in \{1, \dots, n_c\}$ ,  $s \in \{1, \dots, n_s\}$ , and independent permutations  $(\tau^s)$  from  $S_{n_p}$  uniformly at random for  $s \in \{1, \dots, n_s\}$ . With these permutations fixed, for each  $\theta \in \Theta$ , one lets

$$X_i^{\ell,s} := \begin{cases} 1 & \text{if } \sigma^{\ell,s}(i) \leq \zeta(\lambda^\ell n_I), \\ 0 & \text{otherwise,} \end{cases} \quad \text{for } \ell \in \{1, \dots, n_c\}, s \in \{1, \dots, n_s\}, i \in I;$$

$$\xi_{ij}^s := \begin{cases} 1 & \text{if } \tau^s((i,j)) \leq \zeta(\mu n_p), \\ 0 & \text{otherwise,} \end{cases} \quad \text{for } s \in \{1, \dots, n_s\}, (i,j) \in I_2.$$

Minimisation over the parameter space  $\Theta$  is conducted with the MATLAB routine `fminsearchbnd` [7] for constrained optimisation.  $\lambda_{\max}^\ell := \bar{Y}^\ell$  is certainly an upper bound on what  $\lambda^\ell$  any point estimator might estimate ( $\ell \in \{1, \dots, n_c\}$ ) as this is the moment estimate in case  $\mu = 0$ . The upper bound  $\mu_{\max}$  on  $\mu$  is left to the user's judgement.

The last user input in addition to  $n_s$ , Method and  $\mu_{\max}$  is  $n_{\text{opt}}$  which specifies how many different initial states to try in the optimisation runs. We expect an inverse relationship between seeding rates and the contamination rate, given the data. Thus the initial parameter values for  $k \in \{1, \dots, n_{\text{opt}}\}$  are chosen as

$$\lambda_{\text{initial}}^\ell(k) = \left(1 - \frac{k-1}{n_{\text{opt}}}\right) \lambda_{\max}^\ell, \quad \text{for } \ell \in \{1, \dots, n_c\},$$

$$\mu_{\text{initial}}(k) = \chi_{\{n_{\text{opt}} > 1\}} \frac{k-1}{n_{\text{opt}} - 1} \mu_{\max}.$$

## 6.2 Results

In order to test the performance of the proposed estimation procedure, we created a number of synthetic datasets with  $n_c = 3$  colours, different sizes and different, known parameter vectors using Method 1. Tables 1–4 report the results of estimating  $\theta_0 = (\lambda^1, \lambda^2, \lambda^3, \mu)$  using different input settings  $(n_s, n_{\text{opt}}, \mu_{\max})$ .

The two estimators, which are based on Methods 1 and 2 of random number generation, are denoted by  $\hat{\theta}_{n_s, n_I}^{(\text{M1})}$  and  $\hat{\theta}_{n_s, n_I}^{(\text{M2})}$ , respectively. We display the relative bias of the estimators in percentage terms:

$$d^{(\text{M1})} = 100 \left| 1 - \hat{\theta}_{n_s, n_I}^{(\text{M1})} / \theta_0 \right|$$

(the operations are coordinatewise), and analogously,  $d^{(\text{M2})}$  for  $\hat{\theta}_{n_s, n_I}^{(\text{M2})}$ . Finally, we let  $\alpha_{\hat{\theta}_{n_s, n_I}^{(\text{M1})}}$  and  $\alpha_{\hat{\theta}_{n_s, n_I}^{(\text{M2})}}$  denote the value of the objective function  $\alpha$  in (18) at  $\hat{\theta}_{n_s, n_I}^{(\text{M1})}$  and  $\hat{\theta}_{n_s, n_I}^{(\text{M2})}$ , respectively.

The computations were conducted on a laptop computer equipped with a 2.8 GHz Intel Core i7-2640M dual-core processor and 8 GB RAM. Although it is clear that the  $n_{\text{opt}}$  parameter searches and for each, the  $n_s$  simulations lend themselves to parallelisation, our implementation does not benefit from this insight. The columns of  $\alpha_{\hat{\theta}_{n_s, n_I}^{(\text{M1})}}$  and  $\alpha_{\hat{\theta}_{n_s, n_I}^{(\text{M2})}}$  display in brackets running times in seconds for completing the parameter estimation procedure. These times are indicative only and their use for comparisons is limited as less demanding other tasks were also running on the computer

simultaneously. As far as we can tell, the parameter estimation ran in RAM without resorting to swap memory on disk.

We found no definitive answer as to whether Method 1 or 2 is preferable. Table 3 suggests Method 2, but Table 4 is as inconclusive as smaller-sized datasets.

Broadly, the relative bias of the estimates becomes smaller as  $n_I$  grows. From  $n_I = 25 \times 25 = 625$  to  $n_I = 500 \times 500 = 250000$ , the relative bias of the  $\mu$  estimate improves from about 35–40% to below 5%. We have also observed that as  $n_I$  grows, there is ever less need to try several initial states because the solutions tend to converge to the same estimator. In our experience, the existence of local optima that necessitate a greater  $n_{\text{opt}}$  were characteristic of the smaller lattice sizes only.

In the smallest dataset, Table 1, one can observe that  $\lambda^1$  and  $\lambda^2$  are consistently overestimated, whereas  $\lambda^3$  and  $\mu$  are underestimated in all six estimations. This turned out to be due to a quirk of the randomly generated dataset. While  $(\lambda^1, \lambda^2, \lambda^3) = (0.1, 0.05, 0.07)$ , in reality, the dataset had

$$\bar{\mathcal{X}} = \frac{1}{n_I} \left( \sum_{i \in I} \mathcal{X}_i^1, \sum_{i \in I} \mathcal{X}_i^2, \sum_{i \in I} \mathcal{X}_i^3 \right) = (0.1072, 0.0528, 0.0640).$$

One can notice that in Table 1,  $\lambda^1$  and  $\lambda^2$  are overestimated to a greater extent than how much  $\lambda^3$  is underestimated. Then the observed systemic underestimation of  $\mu$  is consistent with this in light of the expected inverse relationship between seeding and contamination described at the end of Section 6.1.

In Table 3, where the lattice size  $n_I = 300 \times 300 = 90000$  is most relevant to our practical application in Section 7,  $(n_s, n_{\text{opt}}) \in \{(2, 8), (4, 4), (8, 2)\}$  allow a comparison of different input choices with approximately identical computational cost.  $(n_s, n_{\text{opt}}) = (4, 4)$  and right behind it  $(8, 2)$  proved to be the best choices, beating  $(n_s, n_{\text{opt}}) = (2, 8)$ . Against the expectations,  $(n_s, n_{\text{opt}}) = (5, 5)$  happened to not improve the estimate with input  $(4, 4)$ . On this lattice size,  $\mu$  is estimated to 10% accuracy with 1–2 hours running time. In Table 4, we get better than 5% accuracy on a larger lattice with 7–8 hours running time.

In Tables 1–4, for fixed  $n_I$ ,  $\alpha_{\hat{\theta}_{n_s, n_I}^{(M1)}}$  and  $\alpha_{\hat{\theta}_{n_s, n_I}^{(M2)}}$  tend to decrease for increasing  $n_s$ . This is reassuring, although not a necessity because it is possible that the synthetic dataset is atypical and more simulations (higher  $n_s$ ) do not make it easier to approximate it. Instead, overfitting might yield the lowest  $\alpha$  values.

For further analysis, we introduce two more symbols. One might consider a trivial estimator which assumes no contamination occurring:  $\hat{\mu} = 0$ ,  $\hat{\lambda}^\ell = \bar{\mathcal{Y}}^\ell$ . The corresponding  $\alpha_{\text{triv}}$  denotes a realisation of  $\alpha$  with parameters from this trivial estimator, computed from  $n_s$  simulations with Method 1 or 2.  $\alpha_{\theta_0}$  denotes a realisation of  $\alpha$  with the true parameter  $\theta_0$  and  $n_s$  simulations.

Table 5 compares  $\alpha_{\hat{\theta}_{n_s, n_I}^{(M1)}}$  and  $\alpha_{\hat{\theta}_{n_s, n_I}^{(M2)}}$ ,  $\alpha_{\text{triv}}$  and  $\alpha_{\theta_0}$  for the four computer-generated datasets of Tables 1–4. Except for the smallest case,  $n_I = 625$ ,  $\alpha_{\theta_0}$  is always smaller than  $\alpha_{\text{triv}}$ , as expected. Whereas  $\alpha_{\theta_0}$  decreases with increasing  $n_I$ ,  $\alpha_{\text{triv}}$  stays roughly constant.  $\alpha_{\hat{\theta}_{n_s, n_I}^{(M1)}}$  and  $\alpha_{\hat{\theta}_{n_s, n_I}^{(M2)}}$  decrease only initially as  $n_I$  increases. One would expect them to be between  $\alpha_{\theta_0}$  and  $\alpha_{\text{triv}}$ , which tends to hold for larger lattice sizes. In reality, their value is much lower than  $\alpha_{\theta_0}$ , but the ratio becomes ever less extreme as  $n_I$  grows. This is indicative of initially very strong, but later ever less pronounced overfitting.

To test the behaviour of the objective function  $\alpha_{\theta_0}$  as  $n_I \rightarrow \infty$ , we generated fresh synthetic datasets of different sizes with a common  $\theta_0 = (0.03, 0.04, 0.05, 0.02)$ . Just generating the single dataset of size  $1000 \times 1000$  took 42 seconds. For this exercise, the single datasets were compared to simulations with common simulation count  $n_s = 10$ . Table 6 shows that both  $\alpha_{\theta_0}$  and  $\tilde{\alpha}(\eta) = \eta^T \eta$  converge to zero, although  $\alpha_{\theta_0}$  has larger values because of the normalisation by  $\Omega$  in (19). This is numerical evidence in support of Propositions 3 and 4, even with fixed  $n_s$ .

$n_s$	$n_{\text{opt}}$	$\mu_{\text{max}}$	$\theta_0$	$\hat{\theta}_{n_s, n_I}^{(\text{M1})}$	$d^{(\text{M1})}$	$\alpha_{\hat{\theta}_{n_s, n_I}^{(\text{M1})}}$	$\hat{\theta}_{n_s, n_I}^{(\text{M2})}$	$d^{(\text{M2})}$	$\alpha_{\hat{\theta}_{n_s, n_I}^{(\text{M2})}}$
10	10	0.1	0.1	0.1292	29.16%	0.0124	0.1277	27.74%	0.0215
			0.05	0.0595	19.00%		0.0637	27.39%	
			0.07	0.0611	12.71%		0.0538	23.11%	
			0.06	0.0433	27.86%		0.0376	37.28%	
50	10	0.1	0.1	0.1289	28.92%	0.0128	0.1267	26.75%	0.00875
			0.05	0.0641	28.26%		0.0611	22.18%	
			0.07	0.0603	13.79%		0.0588	15.94%	
			0.06	0.0394	34.37%		0.0422	29.65%	
100	10	0.1	0.1	0.1350	35.01%	0.0108	0.1259	25.87%	0.0107
			0.05	0.0631	26.18%		0.0623	24.53%	
			0.07	0.0613	12.44%		0.0595	15.00%	
			0.06	0.0360	40.00%		0.0403	32.84%	

Table 1: Six estimates for a synthetic dataset with  $n_I = 25 \times 25 = 625$  vertices ( $n_p = 1776$ ) and  $\theta_0 = (0.1, 0.05, 0.07, 0.06)$ . In this synthetic dataset, the relative frequency of the incidence of the three colours in the seeding is  $\bar{\mathcal{X}} = (0.107, 0.0528, 0.064)$ , while in the contamination-impacted observed data, it is  $\bar{\mathcal{Y}} = (0.15, 0.0752, 0.08)$ . The relative frequency of adjacent vertices having an open edge between them is 0.0574.

$n_s$	$n_{\text{opt}}$	$\mu_{\text{max}}$	$\theta_0$	$\hat{\theta}_{n_s, n_I}^{(\text{M1})}$	$d^{(\text{M1})}$	$\alpha_{\hat{\theta}_{n_s, n_I}^{(\text{M1})}}$	$\hat{\theta}_{n_s, n_I}^{(\text{M2})}$	$d^{(\text{M2})}$	$\alpha_{\hat{\theta}_{n_s, n_I}^{(\text{M2})}}$
20	10	0.05	0.07	0.0734	4.79%	1.55e-05	0.0739	5.52%	0.00044
			0.05	0.0508	1.52%		0.0492	1.66%	
			0.04	0.0390	2.49%		0.0392	2.07%	
			0.03	0.0259	13.77%		0.0262	12.79%	
40	10	0.05	0.07	0.0737	5.28%	5.96e-06	0.0733	4.78%	0.000108
			0.05	0.0505	0.93%		0.0500	0.06%	
			0.04	0.0392	2.05%		0.0393	1.82%	
			0.03	0.0253	15.54%		0.0252	16.00%	

Table 2: Four estimates for a synthetic dataset with  $n_I = 100 \times 100 = 10000$  vertices ( $n_p = 29601$ ) and  $\theta_0 = (0.07, 0.05, 0.04, 0.03)$ . In this synthetic dataset, the relative frequency of the incidence of the three colours in the seeding is  $\bar{\mathcal{X}} = (0.0717, 0.0497, 0.0387)$ , while in the contamination-impacted observed data, it is  $\bar{\mathcal{Y}} = (0.085, 0.0585, 0.0455)$ . The relative frequency of adjacent vertices having an open edge between them is 0.0303.

$n_s$	$n_{\text{opt}}$	$\mu_{\text{max}}$	$\theta_0$	$\hat{\theta}_{n_s, n_I}^{(\text{M1})}$	$d^{(\text{M1})}$	$\alpha_{\hat{\theta}_{n_s, n_I}^{(\text{M1})}}$	$\hat{\theta}_{n_s, n_I}^{(\text{M2})}$	$d^{(\text{M2})}$	$\alpha_{\hat{\theta}_{n_s, n_I}^{(\text{M2})}}$
2	8	0.05	0.05	0.0472	5.57%	0.00241	0.0488	2.43%	0.000149
			0.06	0.0572	4.71%		0.0593	1.18%	
			0.03	0.0308	2.71%		0.0298	0.51%	
			0.02	0.0225	12.59%		0.0220	9.86%	
4	4	0.05	0.05	0.0480	4.03%	0.000968	0.0489	2.12%	6.11e−05
			0.06	0.0578	3.59%		0.0589	1.90%	
			0.03	0.0305	1.65%		0.0298	0.83%	
			0.02	0.0223	11.64%		0.0213	6.45%	
8	2	0.05	0.05	0.0483	3.32%	0.000904	0.0481	3.75%	0.000612
			0.06	0.0586	2.35%		0.0586	2.26%	
			0.03	0.0299	0.40%		0.0301	0.45%	
			0.02	0.0223	11.60%		0.0220	10.11%	
5	5	0.05	0.05	0.0481	3.81%	0.000964	0.0485	2.97%	0.000495
			0.06	0.0580	3.29%		0.0588	1.99%	
			0.03	0.0302	0.79%		0.0303	1.04%	
			0.02	0.0221	10.69%		0.0219	9.72%	

Table 3: Eight estimates for a synthetic dataset with  $n_I = 300 \times 300 = 90000$  vertices ( $n_p = 268801$ ) and  $\theta_0 = (0.05, 0.06, 0.03, 0.02)$ . In this synthetic dataset, the relative frequency of the incidence of the three colours in the seeding is  $\bar{\mathcal{X}} = (0.0498, 0.0595, 0.0299)$ , while in the contamination-impacted observed data, it is  $\bar{\mathcal{Y}} = (0.0558, 0.0667, 0.034)$ . The relative frequency of adjacent vertices having an open edge between them is 0.0198.

$n_s$	$n_{\text{opt}}$	$\mu_{\text{max}}$	$\theta_0$	$\hat{\theta}_{n_s, n_I}^{(\text{M1})}$	$d^{(\text{M1})}$	$\alpha_{\hat{\theta}_{n_s, n_I}^{(\text{M1})}}$	$\hat{\theta}_{n_s, n_I}^{(\text{M2})}$	$d^{(\text{M2})}$	$\alpha_{\hat{\theta}_{n_s, n_I}^{(\text{M2})}}$
1	1	0.04	0.03	0.0336	11.86%	0.0243	0.0331	10.47%	0.025
			0.04	0.0444	10.96%		0.0450	12.59%	
			0.05	0.0553	10.70%		0.0555	10.94%	
			0.02	0.0141	29.27%		0.0136	31.91%	
5	5	0.04	0.03	0.0297	1.08%	0.000894	0.0297	1.04%	0.00106
			0.04	0.0394	1.43%		0.0395	1.13%	
			0.05	0.0518	3.57%		0.0519	3.86%	
			0.02	0.0196	1.94%		0.0194	3.01%	

Table 4: Four estimates for a synthetic dataset with  $n_I = 500 \times 500 = 250000$  vertices ( $n_p = 748001$ ) and  $\theta_0 = (0.03, 0.04, 0.05, 0.02)$ . In this synthetic dataset, the relative frequency of the incidence of the three colours in the seeding is  $\bar{\mathcal{X}} = (0.0299, 0.0402, 0.0503)$ , while in the contamination-impacted observed data, it is  $\bar{\mathcal{Y}} = (0.0336, 0.0451, 0.057)$ . The relative frequency of adjacent vertices having an open edge between them is 0.01996.

$n_I$	$n_p$	$n_s$	$n_{\text{opt}}$	$\alpha_{\theta_0}^{(\text{M1})}$	$\alpha_{\theta_0}^{(\text{M2})}$	$\alpha_{\text{triv}}^{(\text{M1})}$	$\alpha_{\text{triv}}^{(\text{M2})}$	$\alpha_{\hat{\theta}_{n_s, n_I}}^{(\text{M1})}$	$\alpha_{\hat{\theta}_{n_s, n_I}}^{(\text{M2})}$
$25 \times 25$	1776	10	10	1.06	1.77	0.6	0.555	0.0124	0.0215
$25 \times 25$	1776	50	10	0.82	0.919	0.61	0.618	0.0128	0.00875
$25 \times 25$	1776	100	10	1.19	1.09	0.597	0.589	0.0108	0.0107
$100 \times 100$	29601	20	10	0.042	0.0454	0.59	0.61	1.55e−05	0.00044
$100 \times 100$	29601	40	10	0.0496	0.0576	0.598	0.601	5.96e−06	0.000108
$300 \times 300$	268801	2	8	0.0179	0.0048	0.646	0.604	0.00241	0.000149
$300 \times 300$	268801	4	4	0.00948	0.0014	0.653	0.624	0.000968	6.11e−05
$300 \times 300$	268801	8	2	0.00792	0.00659	0.654	0.614	0.000904	0.000612
$300 \times 300$	268801	5	5	0.0104	0.00786	0.651	0.622	0.000964	0.000495
$500 \times 500$	748001	1	1	0.0223	0.0119	0.628	0.621	0.0243	0.025
$500 \times 500$	748001	5	5	0.00626	0.00791	0.633	0.624	0.000894	0.00106

Table 5: A comparison of the values of the objective functions for the true value  $\theta_0$ , for the trivial estimator and for the MSM estimator. The four synthetic datasets used are the same as in Tables 1–4.

Size	$n_I$	$n_p$	$n_s$	$\alpha_{\theta_0}^{(\text{M1})}$	$\alpha_{\theta_0}^{(\text{M2})}$	$\tilde{\alpha}_{\theta_0}^{(\text{M1})}$	$\tilde{\alpha}_{\theta_0}^{(\text{M2})}$
$25 \times 25$	625	1776	10	0.202	0.16	6.65e−05	4.21e−05
$100 \times 100$	10000	29601	10	0.0427	0.0277	7.29e−06	2.57e−06
$300 \times 300$	90000	268801	10	0.00624	0.00807	1.68e−06	1.57e−06
$500 \times 500$	250000	748001	10	0.000799	0.0015	3.87e−07	3.97e−07
$707 \times 707$	499849	1496720	10	0.000521	0.000365	5.08e−07	3.99e−07
$1000 \times 1000$	1000000	2996001	10	0.00127	0.00117	9.78e−08	9.04e−08

Table 6: Realisations of the objective function  $\alpha$  for the true parameter value  $\theta_0$  for different synthetic dataset sizes and of the not normalised variant of the objective function  $\tilde{\alpha}(\eta) = \eta^T \eta$ . Here  $\theta_0 = (0.03, 0.04, 0.05, 0.02)$  across fresh synthetic datasets.

## 7 Cross-contamination rate estimation for digital PCR in lab-on-a-chip microfluidic devices

Our motivation for investigating this problem is the need for quality control in parallelised biochemical experiments run in novel, lab-on-a-chip microfluidic devices for applications in basic research, biotechnology, medical diagnostics and rapid vaccine development. Our collaborators Dr Günter Roth and his group (Centre for Biological Systems Analysis [ZBSA], University of Freiburg) develop such microfluidic devices. The central element of their system is a rectangular well plate with 15 mm edge lengths, with more than 100,000 wells of 19 pℓ volume each. The wells on this chip are arranged in a hexagonal tiling pattern (honeycomb lattice).

Whereas the rival microfluidic technology uses an emulsion of water droplets flowing in an oil medium, this array-based setup fixes a spatial structure, allowing the otherwise neglected analysis of cross-contamination between reaction volumes. Our focus is on evaluating an experiment particularly well suited for this purpose, whose results generalise to other experiments conducted in this lab-on-a-chip device.

In the *digital PCR* experiment, a solution of DNA samples is injected onto the well plate, at such a low concentration that most wells receive 0 or 1 DNA molecule (hence the name *digital*). In the particular case, the solution is a mixture of three different DNA species. We call these template molecules *seeds*. The well plate is covered with a lid (a microscope slide) that is pre-coated with covalently bound DNA primers [15]. The well plate together with the lid serve to insulate the reaction volumes from each other. The DNA templates are amplified in each of the wells independently with a polymerase chain reaction (PCR). In more detail, the template molecules hybridise to the surface-bound primers and the PCR elongates these primers to form the complementary strand of the template. In the next heating step, the templates become resolved, whereas the generated complementary DNA strands stay covalently bound to the surface. The single-strand templates will bind to other surface-bound primers and turn them too into complementary strands via polymerisation. The result of the PCR cycles is that the whole glass surface above the well gets covered with immobilised complementary DNA strands. They mirror the spatial arrangement of the initial seed pattern of the wells.

After the PCR, the three complementary DNA species on the slide are identified via three specifically binding fluorescent hybridisation probes (fluorophores) and their presence or absence can be determined by imaging [16]. In the fluorescent image of the slide (Figure 2), we see either black background (where there was no seed), spots in one of the three primary colours indicating a single seed, and sometimes a mixture of two or three primary colours indicating heterogeneous seeding by multiple seeds. Sometimes we also see clusters of one colour, or an unusually high number of mixed colours, indicating cross-contamination between adjacent wells. This happens when the lid is not fitted tightly and during thermal cycling, liquid exchange occurs between reaction volumes around trapped air bubbles and dust particles. In the readout it remains unclear if two neighbours with the same colour (or a single well with a mix of two colours, which has coloured neighbours) were initiated by two seeds or one contaminated the other (Fig. 2, bottom panel).

For cross-contamination rate estimation for this experimental setup it is necessary to define a mathematical model of the physical process. It has to involve the triangular lattice, which is the dual of the hexagonal tiling, and colouring of its vertices. The total numbers of DNA templates of each type  $\ell \in \{1, \dots, n_c\}$  present in the chip are likely well approximated by  $n_c$  discretised normal random variables. We can safely assume that each well receives a Poisson distributed random number of DNA templates of type  $\ell$  because then due to the superposition property, the total number of type  $\ell$  templates in the chip is also Poisson distributed, which is close to a normal distribution. The Bernoulli distributed  $(X_i^\ell)$  used in our model for seeding are really just a proxy to the either

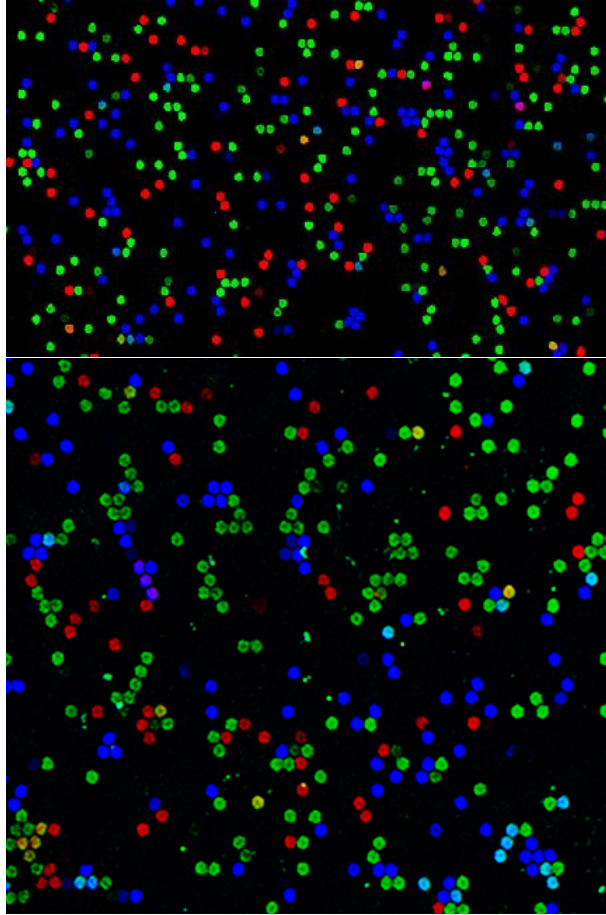


Figure 2: (*top*) Image of a glass slide from a digital PCR experiment with little sign of cross-contamination. [20] (*bottom*) Image of a slide with clustering fluorescent signals and a higher prevalence of cyan and yellow colours, suggesting higher cross-contamination rate.

zero or positive value of the corresponding Poisson distribution. From a value  $\lambda^\ell$  of the Bernoulli parameter, we can infer the parameter  $\tilde{\lambda}^\ell$  of the respective Poisson distribution through the identity  $\lambda^\ell = 1 - e^{-\tilde{\lambda}^\ell}$ .

It is also natural to model the possibility of contamination by open edges. It is a useful shortcut to draw the state of the edges independently of the seeding so that an open edge means only the possibility of propagation, which is contingent on the presence of seeds. There are modelling choices to be made. Contamination might be

- (i) unidirectional (there is the possibility of a pair of independent, oppositely oriented directed edges  $\xi_{i \rightarrow j}$  and  $\xi_{j \rightarrow i}$  between any two adjacent vertices  $i \sim j$ ), or
- (ii) symmetric (undirected edges  $\xi_{ij}$ ).

Open edges might be best represented by

- (1) independent Bernoulli variables, or by
- (2) locally correlated 0–1 random variables.

Contamination might be

- (A) confined to neighbours, or
- (B) it might propagate via a series of open edges.

The choice of (ii,1,B) yields the model put forward in Section 1 (Figure 3). Its strength is that it can use standard percolation theory. Our MSM estimator was developed for this model.

For the quality certification of this lab-on-a-chip device, it is useful to estimate in addition to  $\mu$ , the total number of vertices which belong to a non-trivial component of the percolation graph. These vertices are the wells which were not insulated from their neighbours. Beyond the digital PCR paradigm, in experimental setups where most wells are expected to give some signal, vertices that are connected to any other are likely to give false signals.

An easy upper bound results from noticing that each edge turns at most two additional vertices connected. For small values of  $\mu$ , edges are actually unlikely to share endpoints. The number of edges is distributed according to a binomial distribution with parameters  $n_p$  and  $\mu$ . Therefore the mean number of potentially contaminated vertices can be estimated as

$$\mathbb{E} \left[ \sum_{|C| \geq 2} |C| \right] \leq 2\mu n_p \sim 6\mu n_I$$

where the asymptotic equality holds under the assumption that the boundary of  $I$  is ‘small’. For concrete examples, the conversion from  $n_p$  to  $n_I$  can be accurately determined.

Another approach results by noticing

$$\begin{aligned} \mathbb{E} \left[ \sum_{|C| \geq 2} |C| \right] &= \mathbb{E} \left[ n_I - \sum_{i \in I} \chi_{\{|C(i)|=1\}} \right] \\ &= n_I - n_I(1 - \mu)^6 + e \\ &= \left( 6\mu - 15\mu^2 + \sum_{k=3}^6 \binom{6}{k} (-1)^{k+1} \mu^k \right) n_I + e, \end{aligned}$$

where  $e$  is the correction for boundary vertices.

Simpler cases are given by (i,1,A) and (ii,1,A) where the moments  $\mathbb{E}[Y_i^\ell]$ ,  $\mathbb{E}[Y_i^\ell Y_i^m]$  and  $\mathbb{E}[Y_i^\ell Y_j^\ell]$  ( $\ell \neq m$ ,  $i \sim j$ ) can be computed explicitly. We used MATHEMATICA (Wolfram Research, Inc.) to deal with the many terms, and we report truncations of the complete result for space considerations in the case (i,1,A). It is anticipated in the practical application that  $\mu < \lambda^\ell$  for every  $\ell$ . For non-boundary vertices, under this assumption on the anticipated magnitudes, the dominant terms of the moments of interest in decreasing order are given as

$$\begin{aligned} \mathbb{E}[Y_i^\ell] &= \mathbb{P}(X_i^\ell = 1) + \mathbb{P}(X_i^\ell = 0) \sum_{k=1}^6 \binom{6}{k} \mu^k (1 - \mu)^{6-k} (1 - (1 - \lambda^\ell)^k) \\ &= \lambda^\ell + 6\lambda^\ell \mu - 6(\lambda^\ell)^2 \mu - 15(\lambda^\ell)^2 \mu^2 + \mathcal{O}((\lambda^\ell)^5), \\ \mathbb{E}[Y_i^\ell Y_i^m] &= \mathbb{P}(X_i^\ell X_i^m = 1) + \mathbb{P}(X_i^\ell = 1, X_i^m = 0) \sum_{k=1}^6 \binom{6}{k} \mu^k (1 - \mu)^{6-k} (1 - (1 - \lambda^m)^k) \\ &\quad + \mathbb{P}(X_i^\ell = 0, X_i^m = 1) \sum_{k=1}^6 \binom{6}{k} \mu^k (1 - \mu)^{6-k} (1 - (1 - \lambda^\ell)^k) \\ &\quad + \mathbb{P}(X_i^\ell = X_i^m = 0) \sum_{k=1}^6 \binom{6}{k} \mu^k (1 - \mu)^{6-k} (1 - (1 - \lambda^\ell)^k) (1 - (1 - \lambda^m)^k) \\ &= \lambda^\ell \lambda^m + 18\lambda^\ell \lambda^m \mu - 12((\lambda^\ell)^2 \lambda^m + \lambda^\ell (\lambda^m)^2) \mu + 30\lambda^\ell \lambda^m \mu^2 + \mathcal{O}(\max\{\lambda^\ell, \lambda^m\}^5). \end{aligned}$$

For  $\mathbb{E}[Y_i^\ell Y_j^\ell]$  ( $i \sim j$ ), in the case  $X_i^\ell + X_j^\ell = 1$ , the empty vertex might have been contaminated by the seeded vertex, or it might have been contaminated from its five remaining neighbours. If  $X_i^\ell = X_j^\ell = 0$ , then one can separate cases according to the seeding status of the two shared



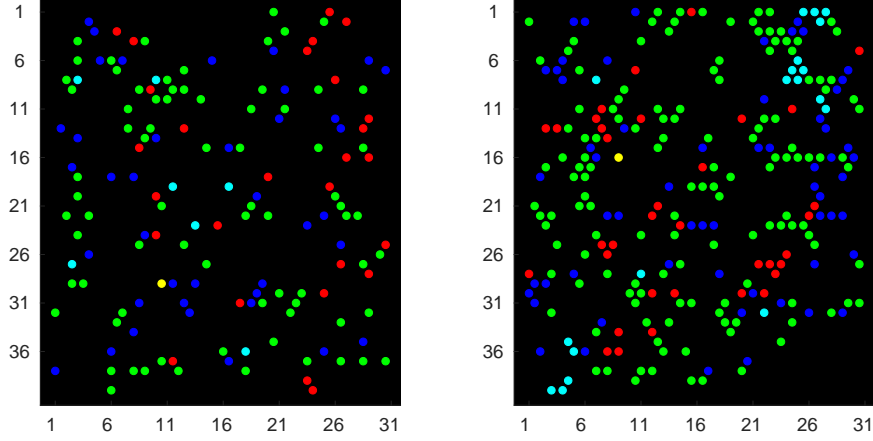


Figure 3: (*left*) Computer simulation of a glass slide from a digital PCR experiment under model (ii,1,B) with  $\theta_0 = (\lambda^{\text{red}}, \lambda^{\text{green}}, \lambda^{\text{blue}}, \mu) = (0.02, 0.07, 0.05, 0.01)$  and relatively little sign of cross-contamination. (*right*) Computer-simulated slide with clustering fluorescent signals and a higher prevalence of cyan colour, suggesting higher cross-contamination rate. Here  $\theta_0 = (0.02, 0.07, 0.05, 0.06)$ .

neighbours of  $i$  and  $j$ . These considerations give

$$\mathbb{E}[Y_i^\ell Y_j^\ell] = (\lambda^\ell)^2 + 2\lambda^\ell \mu + 8(\lambda^\ell)^2 \mu + 2\lambda^\ell \mu^2 - 10(\lambda^\ell)^3 \mu + 9(\lambda^\ell)^2 \mu^2 + \mathcal{O}((\lambda^\ell)^5).$$

These  $n_c^2/2 + 3n_c/2$  moment equations provide the opportunity to estimate the  $n_c + 1$  parameters via the method of moments. Of these, it is  $\mathbb{E}[Y_i^\ell Y_j^\ell]$  where the first term with  $\mu$  is highest up in the magnitude ranking, underpinning the physical intuition that the cooccurrence of a colour in two adjacent vertices is the most informative moment about the contamination rate  $\mu$ .

Notably, the model (ii,1,A) gives exactly the above moment equations if for any  $(i, j) \in I_2$ ,

$$\begin{aligned} \mathbb{P}(\xi_{i \rightarrow j} = 1) &= \mathbb{P}(\xi_{j \rightarrow i} = 1) = \mu && \text{in model (i,1,A), and} \\ \mathbb{P}(\xi_{ij} = 1) &= \mu && \text{in model (ii,1,A).} \end{aligned}$$

The reason is that the propagation of colours is limited to neighbours, so already second neighbours are ruled out. An edge between  $i$  and  $j$  makes a difference in any of the above three moments if and only if  $X_i^\ell + X_j^\ell = 1$ . Say,  $X_j^\ell = 1 = 1 - X_i^\ell$ . Then  $\xi_{j \rightarrow i}$  has the same effect on these moments as  $\xi_{ij}$ , and also the same probability because one can marginalise over the state of  $\xi_{i \rightarrow j}$ . However,  $\mathbb{E}[Y_i^\ell Y_j^\ell Y_i^m Y_j^m]$  would differ between the models (i,1,A) and (ii,1,A). See also the Appendix of [10].

## 8 Discussion and open problems

This paper describes the solution of a statistical problem motivated by a concrete practical need. The mathematical modelling part is solved in one of multiple possible ways, and the choice of (ii,1,B) brings in bond percolation into the statistical model. The percolation is subcritical. The parameter estimation method we propose is the MSM, which gives a point estimate. We prove that it is strongly consistent in the limit as the sample size  $n_I$  tends to infinity. It is an important point that the number of simulations per proposed parameter vector,  $n_s$ , can remain bounded to achieve this result.

What is unusual in our setting is that although the sample size is large, it is not independent (nor identically distributed). Introductory percolation theory is used to upper bound long-range dependencies between the  $n_I$  samples.

We have implemented the method and its accuracy is tested on synthetic datasets in practically relevant parameter ranges. Estimates for wetlab data are to be published by our collaborators Günter Roth and his co-workers in the microfluidics literature.

Parameter estimation in connection with a (static) percolation model is not common in the literature, apart from the quest for the critical value. Dynamic percolation models and dynamic random graphs on a fixed vertex set provide a framework for the contact network in modelling the spread of epidemics. Gilligan and Gibson have been particularly active in studying statistical problems for spatiotemporal models of plant epidemic spread [11, 17]. Gilligan and co-workers also conducted experiments with the fungal pathogen *Rhizoctonia solani* grown in a Petri dish to test how infection probability between a pair of lattice points (that is, the parameter  $\mu$  of percolation in the directed case (i)) depends on their distance and how invasive spread (percolation) probability depends on nutrient availability in lattice points and on the distance between lattice points [2]. They also demonstrated that the random removal (blocking) of sites can hinder and even stop disease spread by driving it subcritical [18].

Beyond the almost sure convergence and the numerical studies with synthetic data, we cannot predict the accuracy of our estimator for instance in terms of confidence intervals. It is known that under regularity conditions, especially that the estimator is continuously differentiable with respect to the parameter  $\theta$ ,  $\sqrt{n_I}(\hat{\theta}_{n_s, n_I} - \theta_0)$  is asymptotically normal with known limiting variance [12, Section 2.3.1]. It is also possible to choose  $\Omega$  optimally, that is, to minimise this asymptotic variance [12, Section 2.3.4]. However, our estimator is not even continuous in  $\theta$  because we use what is called a frequency simulator. It is unknown to us whether it is possible to replace the frequency simulator with some importance sampling to achieve asymptotic normality.

Maximum likelihood estimation (MLE) would have the advantage over MSM that its output is reproducible. Its computational cost might also be lower. Consider the following. We know that black areas have no seeds but we have no information about contamination (edges) in them. We also know that at boundaries between different colours, there is no open edge. Therefore, for a MLE, one needs to establish the probabilities of patches with a fixed colour without knowing which vertices were seeded and which got contaminated only.

We wonder if it is possible by using a generating function that encodes the probabilities of seeding and open edges to compute the total probability that the particular patch was created: each vertex in a patch has been seeded or contaminated from a seed somewhere within the patch. We were only able to derive this generating function for patches that are a linear chain of vertices.

General finite, connected patch shapes (subgraphs) are called (*lattice*) *animals*. Bousquet-Mélou did much work on characterising them via generating functions [5, 6]. Our patches can arise as a disjoint union of adjacent connected components (animals). For our application, it would suffice to develop a recursion which allows one to compute generating functions of small patches (large patches are rare) with a computer algebra system. The difficulty is that the problem is two dimensional,

and a patch must be split in all possible ways into two disjoint parts in the recursion. Any newly added vertex might have been seeded, or contaminated from the rest of the patch, but it might have itself contaminated other empty vertices of the patch.

Notably, the MSM estimator can be turned into an Approximate Bayesian computation (ABC) method very easily. One needs to fix a prior distribution on  $\Theta$  and a small  $\varepsilon > 0$ . The ABC rejection algorithm draws finitely many independent  $\theta \in \Theta$  parameter values from the prior distribution. The objective function (18) is evaluated for each proposed  $\theta$ . The simulations used for the evaluation should no longer use common random numbers but independent ones, and  $n_s$  can be set to one. If the value of the objective function is less than  $\varepsilon$ , then the proposed  $\theta$  is accepted, otherwise it is rejected. This way the set of accepted  $\theta$  is a good approximation of the posterior distribution.

We have not yet tested model fit due to the lack of experimental data. As contamination is caused by the imperfect fit of the glass lid and trapped bubbles and dust, we anticipate that locally positively correlated open edges might be needed in the model. That is, case (ii,2,B) deserves close attention. One way of modelling positive correlations is to apply the Ising model to the edges. Let  $\tilde{\xi}_{ij} = 2\xi_{ij} - 1 \in \{-1, +1\}$ . Then the energy or the Hamiltonian function of a configuration  $\xi$  of open edges is

$$H(\xi) = -J \sum_{i < j < k} (\tilde{\xi}_{ij}\tilde{\xi}_{ik} + \tilde{\xi}_{ij}\tilde{\xi}_{jk} + \tilde{\xi}_{ik}\tilde{\xi}_{jk}) - \tilde{\mu} \sum_{(i,j) \in I_2} \tilde{\xi}_{ij}$$

for some  $J > 0$  and  $\tilde{\mu} < 0$ , and in the first sum, out of the three terms those are missing where an adjacency condition is not met:  $\tilde{\xi}_{ij} = 0$  if  $i \sim j$ , so that every pair of incident edges appears once. The probability of the system being in state  $\xi$  is proportional to  $e^{-\beta H(\xi)}$  for some  $\beta > 0$ . Although we have two new parameters  $J$  and the inverse temperature  $\beta$  in addition to  $\tilde{\mu}$ , the increase in degrees of freedom is really just one,  $\beta J$  and  $\beta \tilde{\mu}$  relative to  $\mu$ .

## 9 Acknowledgements

The authors are grateful to Günter Roth and Christin Rath (ZBSA, University of Freiburg) for proposing the problem, for their relentless help in clarifying details of the experimental protocol and for providing sample images. The authors also thank Robin Ryder (Paris Dauphine University) for suggesting the method of simulated moments and Ed Crane (University of Bristol) and Peter Pfaffelhuber (University of Freiburg) for insights. B. M. thanks the AXA Research Fund for their financial support in the form of a postdoctoral fellowship, and the Isaac Newton Institute for Mathematical Sciences (Cambridge, UK) for support and hospitality during the programme *Stochastic dynamical systems in biology: numerical methods and applications* when work on this paper was undertaken. This work was thereby supported by EPSRC Grant Number EP/K032208/1.

## A Identifiability and numerical estimates of the selected moments

We outline why we conjecture that the parameter  $\theta = (\lambda^1, \dots, \lambda^{n_c}, \mu)$  is identifiable from the moments  $((EY_i^\ell)_{\ell \in \{1, 2, \dots, n_c\}}, (E[Y_i^\ell Y_j^\ell])_{\ell \in \{1, 2, \dots, n_c\}})$  ( $i \sim j$ ). If we focus on just one colour  $\ell$ , then the graph of the function  $(\lambda^\ell, \mu) \mapsto EY_i^\ell$  on the domain  $[0, 1] \times [0, p_c]$  has level curves which go from high  $\lambda^\ell$  and low  $\mu$  to low  $\lambda^\ell$  and high  $\mu$ . In words, the density  $EY_i^\ell$  of colour  $\ell$  is constant if we compensate for a decreasing seeding rate  $\lambda^\ell$  by an appropriately increasing contamination rate  $\mu$ . The function  $(\lambda^\ell, \mu) \mapsto E[Y_i^\ell Y_j^\ell]$  ( $i \sim j$ ) has level curves with the same property.

However, we conjecture that the level curves of  $EY_i^\ell$  and  $E[Y_i^\ell Y_j^\ell]$  do not coincide, instead they intersect. While either one of the two moments narrows down the possible value of the parameter

vector to one of its level curves, the two moments jointly specify the intersection point of two level curves, which uniquely identifies the parameter value  $(\lambda^\ell, \mu)$ .

We provide numerical evidence to back up this claim. For  $n_c = 1$ , we sampled  $EY_i$  and  $E[Y_i Y_j]$  in 142 logarithmically spaced parameter vectors. We made an exception to the logarithmic rule to additionally sample along the line of critical  $\mu$  (Figure A.1). Dataset A contains a broader coverage of 100 parameter vectors. For each of these, we generated independently  $n_s = 5$  realisations of the process on a lattice  $I'$  of size  $300 \times 300$ , and took its central  $100 \times 100$  sublattice  $I \subset I'$  as our data.  $EY_i$  and  $E[Y_i Y_j]$  are estimated as averages over the central sublattice over  $n_s = 5$  realisations.

In Dataset B, 56 parameter vectors are considered which have lower  $\lambda$  values in comparison with Dataset A, save for an overlap of 14 parameter vectors. For each vector, we generated independently  $n_s = 5$  realisations of the process on a lattice  $I'$  of size  $1500 \times 1500$ , and its central  $1000 \times 1000$  sublattice  $I \subset I'$  serves as our data.

The sublattice sizes were selected such that in both datasets, the mean number of seeds is at least 5 in the central sublattice used for sampling, even for their respective lowest  $\lambda$  values ( $\lambda = 5 \times 10^{-4}$  in Dataset A, and approximately  $5.23 \times 10^{-6}$  in Dataset B). At the larger lattice size used for Dataset B, for  $\mu$  values larger than what we tested, the step of finding the connected open components to generate the data became prohibitively time consuming.

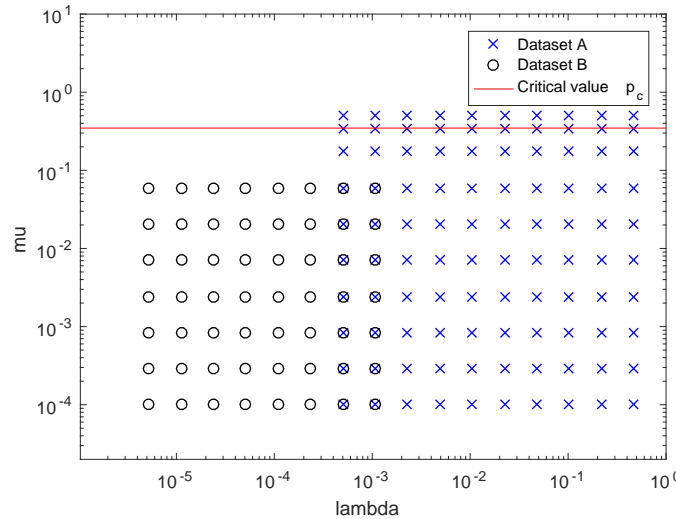


Figure A.1: Sampled parameter values  $\theta = (\lambda, \mu)$ . Dataset A spans  $[5 \times 10^{-4}, 0.4676] \times [10^{-4}, 0.5]$  and Dataset B spans  $[5.23 \times 10^{-6}, 0.00107] \times [10^{-4}, 0.0595]$ .

Figures A.2–A.4 display graphs and level curves of the two coordinates of

$$(\lambda, \mu) \mapsto \left( \frac{1}{n_s n_I} \sum_{s=1}^{n_s} \sum_{i \in I} Y_i^s, \frac{1}{n_s n_p} \sum_{s=1}^{n_s} \sum_{(i,j) \in I_2} Y_i^s Y_j^s \right).$$

Close observation of the level curves seems to show that those in Figure A.2 fan out with different slopes from a smaller region, while those in Figure A.3 are closer to parallel. This supports our conjecture that level curves of one type intersect level curves of the other type in exactly one point, giving identifiability, except perhaps for a null set or otherwise small subset of  $\Theta$  where the two types of level curves coincide.

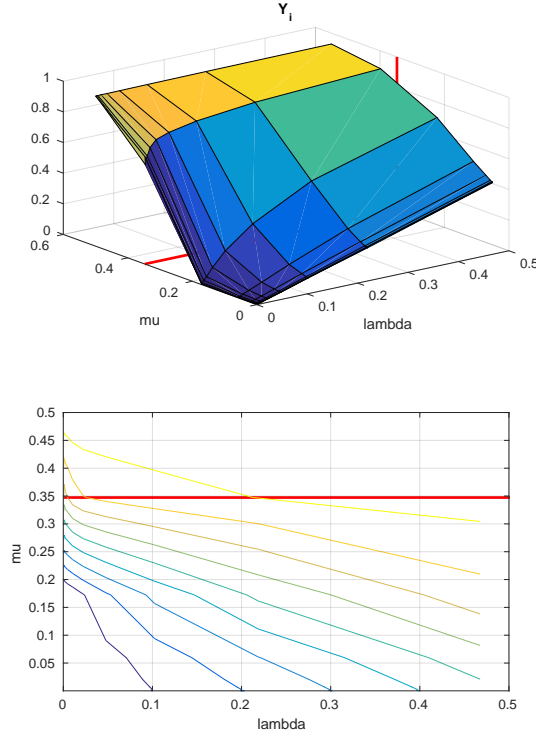


Figure A.2: (*top*) Empirical means of  $Y_i$  for the various parameter vectors of Dataset A. (*bottom*) Level curves of this function. The red lines mark the critical value  $p_c$ .

## References

- [1] Michael Aizenman and Charles M. Newman. Tree graph inequalities and critical behavior in percolation models. *Journal of Statistical Physics*, 36(1–2):107–143, 1984. doi: 10.1007/BF01015729.
- [2] Douglas J. Bailey, Wilfred Otten, and Christopher A. Gilligan. Saprotrophic invasion by the soil-borne fungal plant pathogen *Rhizoctonia solani* and percolation thresholds. *New Phytologist*, 146(3):535–544, June 2000. ISSN 1469-8137. doi: 10.1046/j.1469-8137.2000.00660.x. URL <http://dx.doi.org/10.1046/j.1469-8137.2000.00660.x>.
- [3] Felix Beck. Parameter estimation in a percolation model with coloring. Master’s thesis, Institute for Mathematics, University of Freiburg, Germany, 2015.
- [4] J. van den Berg and H. Kesten. Inequalities with applications to percolation and reliability. *Journal of Applied Probability*, 22(3):556–569, Sep 1985. ISSN 00219002. doi: 10.2307/3213860. URL <http://www.jstor.org/stable/3213860>.
- [5] Mireille Bousquet-Mélou. New enumerative results on two-dimensional directed animals. *Discrete Mathematics*, 180:73–106, 1998. doi: 10.1016/S0012-365X(97)00109-X.
- [6] Mireille Bousquet-Mélou and Andrew Rechnitzer. Lattice animals and heaps of dimers. *Discrete Mathematics*, 258:235–274, 2002. doi: 10.1016/S0012-365X(02)00352-7.

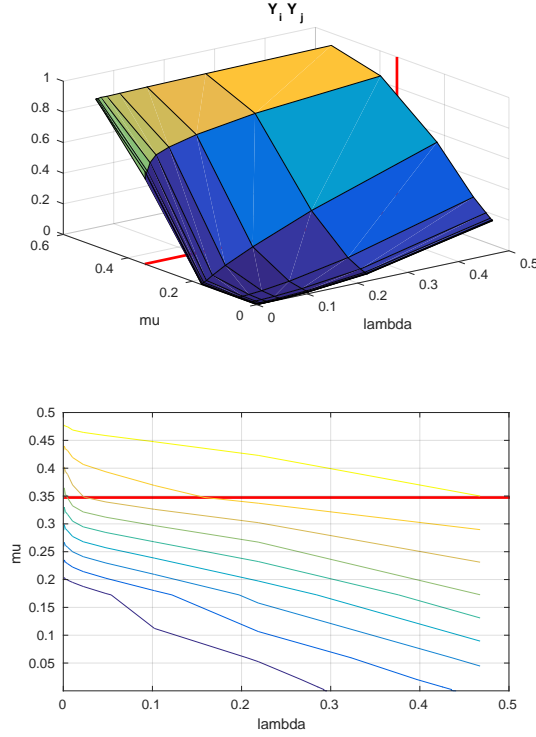


Figure A.3: (*top*) Empirical means of  $Y_i Y_j$  ( $i \sim j$ ) for the various parameter vectors of Dataset A. (*bottom*) Level curves of this function. The red lines mark the critical value  $p_c$ .

- [7] John D’Errico. `fminsearchbnd`, `fminsearchcon` MATLAB files, 2012. URL <http://uk.mathworks.com/matlabcentral/fileexchange/8277-fminsearchbnd--fminsearchcon>.
- [8] Nasrollah Etemadi. On the laws of large numbers for nonnegative random variables. *Journal of Multivariate Analysis*, 13(1):187–193, 1983. ISSN 0047-259X. doi: 10.1016/0047-259X(83)90013-1. URL <http://www.sciencedirect.com/science/article/pii/0047259X83900131>.
- [9] C. M. Fortuin, P. W. Kasteleyn, and J. Ginibre. Correlation inequalities on some partially ordered sets. *Communications in Mathematical Physics*, 22(2):89–103, 1971. ISSN 0010-3616. doi: 10.1007/BF01651330.
- [10] H. L. Frisch and J. M. Hammersley. Percolation processes and related topics. *Journal of the Society for Industrial and Applied Mathematics*, 11(4):894–918, 1963. ISSN 03684245. URL <http://www.jstor.org/stable/2946482>.
- [11] G. J. Gibson, W. Otten, J. A. N. Filipe, A. Cook, G. Marion, and C. A. Gilligan. Bayesian estimation for percolation models of disease spread in plant populations. *Statistics and Computing*, 16(4):391–402, 2006. ISSN 1573-1375. doi: 10.1007/s11222-006-0019-z. URL <http://dx.doi.org/10.1007/s11222-006-0019-z>.
- [12] Christian Gouriéroux and Alain Monfort. *Simulation-based econometric methods*. Oxford University Press, Oxford, UK, 2002.

- [13] Christian Gouriéroux and Alain Monfort. Simulation based inference in models with heterogeneity. *Annales d'Économie et de Statistique*, 20–21:69–107, 1991. ISSN 0769489X. URL <http://www.jstor.org/stable/20075807>.
- [14] Geoffrey Grimmett. *Percolation*. Grundlehren der mathematischen Wissenschaften. Springer Berlin Heidelberg, 1999.
- [15] Jochen Hoffmann, Sebastian Hin, Felix von Stetten, Roland Zengerle, and Günter Roth. Universal protocol for grafting PCR primers onto various lab-on-a-chip substrates for solid-phase PCR. *RSC Advances*, 2:3885–3889, 2012. doi: 10.1039/c2ra01250b.
- [16] Jochen Hoffmann, Martin Trotter, Felix von Stetten, Roland Zengerle, and Günter Roth. Solid-phase PCR in a picowell array for immobilizing and arraying 100 000 PCR products to a microscope slide. *Lab on a Chip*, 12:3049–3054, 2012. doi: 10.1039/c2lc40534b.
- [17] Jonathan J. Ludlam, Gavin J. Gibson, Wilfred Otten, and Christopher A. Gilligan. Applications of percolation theory to fungal spread with synergy. *Journal of The Royal Society Interface*, 9(70):949–956, 2012. ISSN 1742-5689. doi: 10.1098/rsif.2011.0506. URL <http://rsif.royalsocietypublishing.org/content/9/70/949>.
- [18] Wilfred Otten, Douglas J. Bailey, and Christopher A. Gilligan. Empirical evidence of spatial thresholds to control invasion of fungal parasites and saprotrophs. *New Phytologist*, 163(1): 125–132, July 2004. ISSN 1469-8137. doi: 10.1111/j.1469-8137.2004.01086.x. URL <http://dx.doi.org/10.1111/j.1469-8137.2004.01086.x>.
- [19] David Pollard. *Convergence of stochastic processes*. Springer Series in Statistics. Springer, New York, 1984. ISBN 978-1-4612-9758-1. doi: 10.1007/978-1-4612-5254-2.
- [20] Christin Rath. DNA-Kopierprozess mit Thrombin-Aptamer Mikroarrays (in German). Master’s thesis, Centre for Biological Systems Analysis (ZBSA), Faculty of Biology, University of Freiburg, Germany, 2014.
- [21] M. F. Sykes and J. W. Essam. Exact critical percolation probabilities for site and bond problems in two dimensions. *Journal of Mathematical Physics*, 5(8):1117–1127, 1964. doi: 10.1063/1.1704215. URL <http://scitation.aip.org/content/aip/journal/jmp/5/8/10.1063/1.1704215>.
- [22] Sara van de Geer. *Applications of empirical process theory (Empirical processes in M-estimation)*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, Cambridge, UK, 2010.

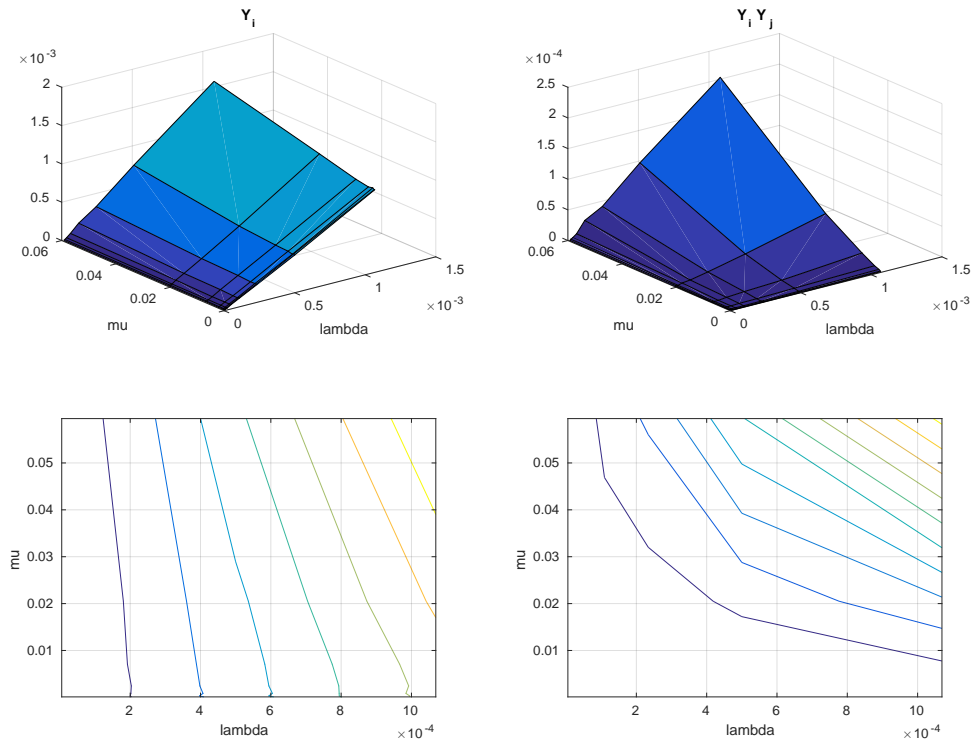


Figure A.4: (*top left*) Empirical means of  $Y_i$  for the various parameter vectors of Dataset B. (*bottom left*) Level curves of this function. (*top right*) Empirical means of  $Y_i Y_j$  ( $i \sim j$ ) for the various parameter vectors of Dataset B. (*bottom right*) Level curves of this function.